# Multi-Agent Reinforcement Learning for Resources Allocation Optimization: A Survey

Mohamad A. Hady[1], Siyi Hu[1*], Mahardhika Pratama[1],
Jimmy Cao[1], Ryszard Kowalczyk[1]

[1]STEM, University of South Australia, Mawson Lakes Blvd, Mawson Lakes, 5095, South Australia, Australia.

*Corresponding author(s). E-mail(s): Siyi.Hu@unisa.edu.au;
Contributing authors: mohamad.hady@mymail.unisa.edu.au;
Dhika.Pratama@unisa.edu.au; jimmy.cao@unisa.edu.au;
Ryszard.Kowalczyk@unisa.edu.au;

**Abstract**

Multi-Agent Reinforcement Learning (MARL) has become a powerful framework for numerous real-world applications, modeling distributed decision-making and learning from interactions with complex environments. Resource Allocation Optimization (RAO) benefits significantly from MARL's ability to tackle dynamic and decentralized contexts. MARL-based approaches are increasingly applied to RAO challenges across sectors playing pivotal roles to Industry 4.0 developments. This survey provides a comprehensive review of recent MARL algorithms for RAO, encompassing core concepts, classifications, and a structured taxonomy. By outlining the current research landscape and identifying primary challenges and future directions, this survey aims to support researchers and practitioners in leveraging MARL's potential to advance resource allocation solutions.

**Keywords:** Multi-Agent Reinforcement Learning, Resource Allocation Optimization

# 1 Introduction

Multi-agent reinforcement learning (MARL) has quickly become an essential area of research, providing effective solutions for distributed decision-making in dynamic and decentralized environments. As multiple agents interact and learn in shared settings, MARL addresses the complexities of real-world applications, especially in situations
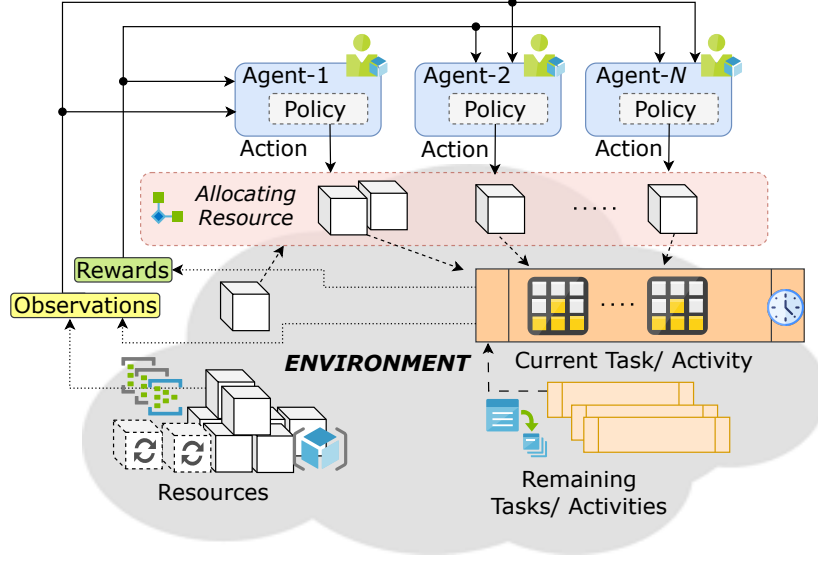
**Fig. 1**: MARL solution for RAO. Any resources can be allocated by several agents to complete tasks or activities. Each agents has its own policy to handle resource allocation determined by the rewards and observations of the whole system states that may come from resources and tasks situations.

with non-stationary and evolving conditions Ning and Xie (2024); Hao et al. (2023); Nguyen et al. (2020). In parallel, Resource Allocation Optimization (RAO) has gained significant attention, as optimizing resource distribution—such as time, energy, network bandwidth, and computational power—can enhance efficiency and effectiveness across a variety of fields Wei et al. (2021); Liu et al. (2023); Allahham et al. (2022).

MARL is particularly suited to tackling RAO challenges, as it enables decentralized, adaptive decision-making. This ability is critical in industries like telecommunications, energy management, cloud computing, and transportation, where efficient resource management plays a vital role Wong et al. (2023); Zabihi et al. (2023). For example, in cloud computing, MARL-based algorithms can optimize resource scheduling and load balancing, leading to improved system performance and reduced costs. Similarly, global supply chains benefit from efficient resource allocation, boosting productivity and reducing expenditures Jiang and Sheng (2009); Ren et al. (2022). In power grids, where renewable energy sources are increasingly integrated, MARL enables dynamic energy resource allocation to balance supply and demand Zhang et al. (2023). Transportation networks also leverage MARL's adaptive capabilities, with applications such as traffic signal control in cities to alleviate congestion and reduce emissions Wu et al. (2020). A generic illustration of MARL solution for RAO framework is provided in Fig. 1 that describes how any resources are allocated to a certain task or activity with multi-agent decision makers.
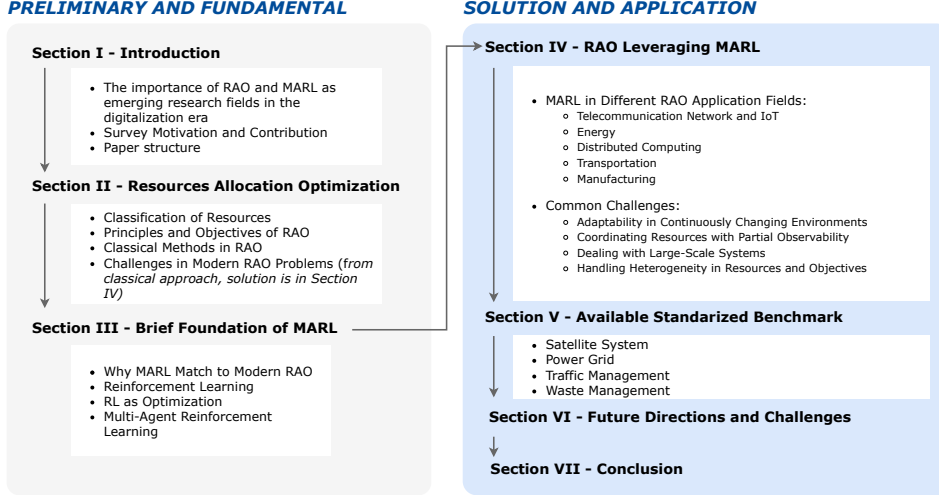
**Fig. 2**: Complete manuscript body structure used in this survey: Preliminary section mainly covers the fundamental of RL and MARL methods. Then, we introduce the concept of RAO, its classical solution and highlight the challenges in the recent trends which can be solved using MARL algorithm and discussed in the RAO leveraging MARL section.

Historically, RAO problems have been addressed through classical optimization and heuristic methods Halabian (2019a). However, these methods often lack of the flexibility and scalability required in complex, real-time environments Sarah et al. (2023). Reinforcement learning (RL), and particularly MARL, has emerged as a powerful alternative. While RL has proven effective for single-agent optimization problems, MARL extends this to systems with multiple interacting agents, each learning and adapting in real time. This shift from single-agent RL to MARL enables more scalable, decentralized approaches suitable for the increasing complexity of modern RAO scenarios Lei et al. (2020); Noor-A-Rahim et al. (2020); Chen et al. (2021).

Despite rising interest in MARL and its RAO applications, there is a lack of comprehensive surveys that focus specifically on this intersection. Existing surveys cover related topics—such as RL in energy systems Yu et al. (2021), task allocation in multi-robot systems Orr and Dutta (2023), and resource management in wireless networks Feriani and Hossain (2021)—but do not offer an extensive review of MARL-driven RAO across different industries. This survey fills that gap by providing a focused review on the use of MARL in RAO, with the following primary contributions:

- Mapping the current landscape of MARL algorithms and frameworks used in RAO, providing researchers with a consolidated resource.
- Offering a systematic review to synthesize advancements, highlight trends, and identify challenges and opportunities unique to MARL applications in RAO.
- Categorizing recent literature in MARL training frameworks and their application areas.

3

- Listing real-world benchmarks and available testbeds for RL and MARL algorithm development in RAO.

This survey is both timely and necessary, serving as a foundational reference for researchers and practitioners in the field. By focusing on recent developments, selecting high-impact studies, and examining critical aspects such as non-stationary, scalability, agent communication, and coordination, we aim to provide a comprehensive overview of current MARL research in the context of RAO. The remainder of this survey is organized as follows (see Fig. 2): *Section 2* categorizes types of resources and reviews classical RAO approaches, including linear programming, heuristic optimization, and game theory. It highlights their applications and limitations in dynamic, large-scale, and decentralized systems. *Section 3* introduces the fundamentals of reinforcement learning (RL) and extends them to multi-agent reinforcement learning (MARL), with a focus on their applicability to decentralized and dynamic RAO scenarios. *Section 4* surveys MARL solutions for RAO-related applications and challenges. It reviews the use of MARL in domains such as telecommunications, energy systems, distributed computing, transportation, and manufacturing, and discusses challenges including dynamic environments, partial observability, scalability, and resource heterogeneity. This section also examines key frameworks such as CTCE, DTDE, and CTDE, along with emerging approaches like graph-based MARL. *Section 5* introduces representative benchmarks used to evaluate MARL performance in RAO settings, including satellite missions, power grids, container management, and traffic systems. *Section 6* outlines future research directions, emphasizing the need for improved scalability, real-time adaptability, and agent coordination. Potential solutions include hierarchical MARL and mean-field approximation techniques.

# 2 Resource Allocation Optimization

Resource Allocation Optimization (RAO) is a significant area of research across many fields, focusing on the challenge of distributing resources among agents or tasks to improve system efficiency, productivity, or fairness. Resource allocation is the structured distribution of finite resources among tasks or activities to meet a defined objective, such as efficiency or cost. This process is fundamental across various fields, including energy management, cloud computing, manufacturing, and telecommunications, etc. The effective allocation requires identification of available resources, assessment of task demands, and allocation of resources in a way that aligns with system goals, often under constraints such as time, budget, or capacity Feriani and Hossain (2021). The primary objective of RAO is to allocate resources as effectively as possible, balancing competing demands to optimize specific goals, such as minimizing delays, maximizing throughput, reducing costs, or ensuring fairness among users Tang et al. (2015). In industries like telecommunications, cloud computing, energy distribution, manufacturing, and transportation, efficient allocation of resources such as bandwidth, computational power, energy, and physical assets is crucial to maintaining high performance while minimizing costs and reducing waste. As systems become more complex and interconnected with advancements like the Internet of Things (IoT)

and digitalization in large-scale distributed systems, the demand for effective resource allocation becomes even more critical.

## 2.1 Principles and Objectives of RAO

Resource location introduces additional complexity in allocation strategies. Non-distributed resources are concentrated in a single location or managed under a centralized authority. While these systems minimize communication overhead and latency, they often face scalability challenges and are susceptible to single points of failure. Conversely, distributed resources are spread across multiple locations or nodes and require coordination across interconnected systems. This setup offers scalability, fault tolerance, and flexibility, making distributed systems ideal for dynamic environments such as cloud computing, multi-robot coordination, and smart grids. However, distributed systems also pose challenges like communication overhead and coordination latency, necessitating advanced algorithms for efficient resource management Zhang and Debroy (2023); Huang et al. (2023).

### 2.1.1 Resource Allocation Process

In practice, resource allocation strategies are guided by specific characteristics of the resources themselves. For instance, in distributed cloud computing, resources like processing power are spread across multiple servers, necessitating allocation methods that balance load across the network to enhance response times and availability Jiang (2015). In contrast, centralized resources in manufacturing, such as machinery, are managed within a single facility to minimize downtime and improve throughput. For divisible resources, such as bandwidth in telecommunications, allocations can vary to meet user demand, reducing latency and enhancing throughput. Indivisible resources, such as individual machines or personnel, require discrete allocation, where entire units are assigned based on task requirements. Finally, renewable resources, like energy from solar panels, are allocated cyclically to maintain availability without depletion, while non-renewable resources, such as a fixed budget, require careful allocation to support long-term goals Wang et al. (2021).

***Problem Formulation:***
Consider a set of tasks indexed by $i$, where $i = 1, 2, \ldots, n$. Let $x_i$ denote the amount of resource allocated to the $i$-th task. The total amount of resource available, $N$, constrains the allocation as: $\sum_{i=1}^{n} x_i \leq N$. Additionally, task-specific limits on resource allocation, represented by lower and upper bounds $l_i$ and $u_i$, are applied as: $l_i \leq x_i \leq u_i, i = 1, 2, \ldots, n$. These constraints can be incorporated directly into the optimization model, allowing for efficient adjustments according to operational needs.

### 2.1.2 Objective of RAO

The objective function for RAO, denoted $f(x_1, x_2, \ldots, x_n)$, is formulated to achieve optimal resource distribution by minimizing costs or maximizing benefits. This can be represented as follows Ibaraki and Katoh (1988); Ushakov (2013):

$$\text{maximizing } f(x_1, x_2, \ldots, x_n)$$

$$\text{subject to} \quad \sum_{i=1}^{n} x_i \leq N, \quad l_i \leq x_i \leq u_i, \; i = 1, 2, \ldots, n, \tag{1}$$

where $N$ represents the total available amount of the resource. In cases where the objective is to maximize profit, the problem can be re-framed by minimizing $-f$, as maximizing $f$ is equivalent to minimizing its negative.

The structure of the objective function $f(x_1, x_2, \ldots, x_n)$ is often customized to fit the specific requirements of the application, and it may take various forms Patriksson (2008):

- Separable function: A common structure where the objective is expressed as the sum of individual cost functions for each resource, such as $\sum_{i=1}^{n} f_i(x_i)$.
- Convex function: When each $f_i$ is convex, enabling the use of convex optimization techniques for efficiency.
- Minimax or maximin objectives: In some scenarios, the goal is to minimize the maximum cost, $\max_i f_i(x_i)$, or maximize the minimum cost, $\min_i f_i(x_i)$, which can help balance resource allocation across tasks.
- Fairness-oriented function: A fairness-focused allocation can be achieved by minimizing a function $g(\max_i f_i(x_i),$
  $\min_i f_i(x_i))$, where $g$ is a non-decreasing function, thereby balancing extremes in resource distribution.

## 2.2 Resources Properties

In RAO, the properties of resources play a pivotal role in determining effective allocation strategies. Resources can be broadly classified by their *divisibility*, *duration*, and *location*, each influencing management approaches and system design. These classifications are illustrated in Fig. 3. Discrete resources—such as physical machines, servers, or personnel—are indivisible and allocated as whole units. These are typically modeled using Integer Programming (IP), where the resource variable $x_j$ must take non-negative integer values ($x_j \in \mathbb{Z} \geq 0$). In contrast, continuous resources—such as bandwidth, processing power, or memory—can be allocated in fractional amounts. These are modeled with real-valued variables ($x_j \in \mathbb{R} \geq 0$), enabling more flexible and fine-grained allocation strategies.

The duration of resource availability significantly impacts allocation methods, as resources may be either renewable or non-renewable. Renewable resources—such as solar energy or CPU time—replenish over time and require time-dependent constraints for effective management ($x_i(t) \leq R_i(t), \; \forall t$). In contrast, non-renewable resources—such as budgets, storage, or fossil fuels—are finite and subject to cumulative upper-bound constraints ($\sum_t x_i(t) \leq R_i$). These distinctions are essential for modeling real-world systems, where renewable and non-renewable resources often coexist and must be managed in a coordinated and efficient manner. Recognizing the temporal nature of resource availability is therefore critical for designing appropriate allocation strategies.
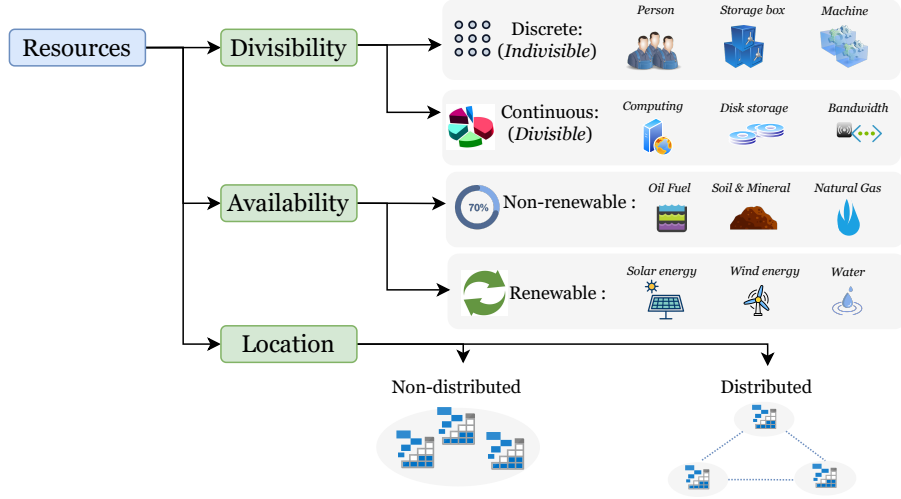
**Fig. 3**: Resource classification by divisibility, availability, and location properties.

Effectively managing resource allocation requires a structured approach to align resources with the specific demands of a system. This process involves understanding the type of resources, whether distributed or non-distributed and the needs of the agents or processes involved. Distributed resources, which are common in cloud computing, allow for flexible and scalable allocation across users and systems, though they introduce challenges in managing latency and coordination Halabian (2019b); Sadat-diynov et al. (2023). In contrast, non-distributed resources, typical in manufacturing, often rely on a centralized allocation framework that provides direct oversight but may experience bottlenecks as demands increase.

Given these distinctions, the allocation strategy should be adapted to the resource type to maintain system efficiency. This approach requires assessing availability, demand, and constraints to ensure that resources are distributed effectively to meet system objectives, such as maximizing throughput or minimizing idle time. The following section outlines a general process for resource allocation, demonstrating how these principles apply to real-world scenarios. Fig. 3 illustrates the distribution of resources in different properties.

### 2.3 Classical Methods in RAO

This section describes several classical methods to understand their limitations and potential applications in RAO (see the full list of methods in Table 1).

#### 2.3.1 Linear Programming

Linear Programming (LP) is one of the foundational techniques in optimization, particularly effective for RAO problems with linear relationships between variables. LP optimizes a linear objective function under linear equality and inequality constraints.

**Table 1**: Summary of classical approaches for RAO Solution

| Category | Algorithm | Representative Work |
|---|---|---|
| Linear Programming | Mixed Integer Linear Programming | Saaty et al. (2003) |
| Heuristic Optimization | Simulated Annealing | Spinellis et al. (2000); Suman and Kumar (2006); Attiya and Hamam (2006); Bi et al. (2020); Kosanoglu et al. (2024) |
| | Genetic Algorithms | Alcaraz and Maroto (2001); Cardon et al. (2000); Tseng et al. (2017); Gao et al. (2020); Shao et al. (2024) |
| | Particle Swarm Optimization | Bratton and Kennedy (2007); Gong et al. (2012); Lin and Chiu (2018); Liu et al. (2022) |
| | Fuzzy Logic Based | Xu et al. (2008); Wu et al. (2018); Khan et al. (2019); Zhang et al. (2021) |
| Game Theory | Cooperative Game | Khan and Ahmad (2006); Zhang et al. (2015) |
| | Non-Cooperative Game | Khan and Ahmad (2006); Ye and Chen (2013) |

This makes it useful for traditional RAO problems such as cost minimization, profit maximization, and efficient allocation of resources like time, money, or physical assets Saaty et al. (2003).

While LP is suitable for problems with manageable numbers of variables and constraints, modern systems like cloud computing or telecommunications often involve high-dimensional settings with large numbers of constraints and variables, making LP computationally expensive for large-scale RAO in real-time contexts.

LP is inherently a centralized approach, assuming a single entity with full knowledge of resources, constraints, and objectives. However, many RAO scenarios involve decentralized systems, like multi-agent systems in smart grids or IoT networks, where LP's centralized nature can lead to inefficiencies due to communication overhead and the need to gather global information.

### 2.3.2 Heuristic Optimization

Heuristic algorithms provide a flexible alternative for RAO by focusing on finding acceptable, near-optimal solutions within a reasonable time frame. These methods are particularly useful for complex, large-scale, or time-sensitive RAO scenarios where exact solutions are impractical. While heuristics do not guarantee an optimal solution, they offer a practical trade-off between accuracy and computational efficiency, making them suitable for real-time applications.

- *Simulated Annealing (SA)*: Inspired by the annealing process, SA probabilistically explores the solution space, refining the search to find near-optimal solutions, especially useful in complex RAO scenarios Kirkpatrick et al. (1983); Spinellis et al. (2000).

- *Genetic Algorithms (GA)*: GA uses principles of natural selection to evolve solutions, making it suitable for large or complex RAO problems with vast search spaces, such as those found in multi-agent systems Cardon et al. (2000) and cloud resource allocation Tseng et al. (2017).
- *Particle Swarm Optimization (PSO)*: PSO simulates social behaviors, where particles (potential solutions) adjust based on their experiences and those of their neighbors. This approach is particularly well-suited for distributed optimization problems in RAO Kennedy and Eberhart (1995).
- *Fuzzy Logic-Based Algorithms*: These algorithms handle uncertainty in RAO by using linguistic variables and fuzzy rules, providing flexible allocation even when resource demand and availability are imprecise Xu et al. (2008).

### 2.3.3 Game Theory

Game theory models RAO as a strategic interaction among agents, each with potentially competing resource needs. In non-cooperative RAO scenarios, each agent aims to maximize its utility independently Khan and Ahmad (2006); Ye and Chen (2013). For example, in wireless communication, users share a common spectrum, each aiming to maximize data rates while accounting for interference from others Cesana et al. (2008). The resulting Nash equilibrium represents an optimized allocation where no user can unilaterally improve their outcome, though it may not achieve global optimization.

In cooperative RAO scenarios, agents may form coalitions to share resources in ways that improve system-wide outcomes. Cooperative game theory models these coalitions, with the Shapley value providing a way to distribute gains fairly based on each agent's contribution to the coalition Khan and Ahmad (2006); Zhang et al. (2015). This approach is particularly relevant in decentralized RAO, such as multi-agent systems where autonomous agents either compete or cooperate for shared resources Zhang et al. (2012). By capturing interactions among agents, game theory provides a structured framework for optimizing resource sharing in settings where individual and group goals intersect.

## 2.4 Limitation and Challenge in Classical RAO Approach

Classical RAO approaches, including Linear Programming, heuristic optimization, and game theory, provide essential tools for resource allocation. However, modern RAO problems present challenges that these methods cannot fully address. As systems become increasingly complex, interconnected, and dynamic, classical methods struggle with issues like scalability, adaptability, and decentralization. Key challenges include:

- **Continuous and Rapid Changing Issues**: Modern RAO environments, such as power grids, cloud computing, and telecommunications networks, are characterized by rapidly changing resource demands and availability. Classical methods, which assume static or semi-static conditions, struggle to adapt to these dynamic settings. For example, cloud computing experiences sudden fluctuations in demand for processing power, storage, and bandwidth due to varying user activities or service request spikes Zhang and Debroy (2023), while power grids must balance supply and

demand in real time to mitigate voltage fluctuations through optimal power management Sun and Qiu (2021); Alam et al. (2016). Similarly, next-generation mobile communication networks face challenges in real-time adaptability and dynamic load distribution Wang et al. (2024). Applications such as autonomous vehicles, financial trading, and smart grids require rapid decision-making under strict time constraints, but classical methods, which rely on solving complex equations or iterative processes, often cannot provide near-instant solutions Singh et al. (2017). Furthermore, many modern systems operate under significant uncertainties, such as unpredictable network congestion in telecommunications or changing energy supplies and demands influenced by external factors Suzuki et al. (2022).

- **Decentralization and Partial Observability**: In many modern RAO settings, such as multi-agent networks, smart grids, and IoT systems, resource allocation decisions are increasingly decentralized and distributed, with agents making independent decisions based on local information Halabian (2019a); Liao et al. (2020); Hu et al. (2020, 2024). Classical, centralized approaches are not well-suited for these scenarios, as they assume full system observability and a central entity managing allocation. In decentralized systems, each agent operates with partial knowledge of the environment and must balance its own objectives against those of other agents. Handling partial observability requires advanced, distributed algorithms that can coordinate effectively without relying on global information, reducing the need for extensive communication overhead.

- **Scalability Issues**: As systems grow larger, resource allocation complexity escalates, often exceeding the capabilities of classical optimization techniques. Techniques like linear programming and heuristics become computationally infeasible when applied to large-scale systems, such as cloud computing networks, fog computing, or power grids, where thousands or millions of components (e.g., servers, devices, or users) need to share limited resources Costa et al. (2022); Gao et al. (2022b, 2023). The computational costs of classical methods scale poorly with the number of variables and constraints, which leads to significant delays and inefficiencies. Many classical approaches rely on centralized control, where a single entity manages resource allocation based on global system information. In large-scale systems, centralized decision-making is impractical due to bandwidth, latency, and processing limitations, causing bottlenecks and long delays. Techniques like mean-field approximations have been developed to address these scalability challenges, though they remain limited in handling highly interconnected systems Yang et al. (2018); Wang et al. (2020).

- **Heterogeneity of Resources and Objectives**: Modern RAO problems often involve diverse resource types and multi-objective optimization, such as minimizing cost while maximizing efficiency and ensuring fairness. Unlike classical problems that assume homogeneous resources, modern systems must consider a range of resource types (e.g., bandwidth, energy, memory) with unique characteristics and constraints Gao et al. (2022a); Zhao et al. (2023). Additionally, multi-objective optimization requires balancing competing goals, which classical methods often cannot handle without significant adjustments, making them impractical for complex, multi-dimensional RAO tasks.

**Table 2**: Classical Approaches vs. MARL in RAO

| Aspect | Limitation of Classical Approaches | Advantage of MARL |
|---|---|---|
| Adaptability | Assume static or predictable systems; not reactive to changes in real-time | Agents adapt through continuous learning and interaction with dynamic environments |
| Scalability | Poor performance as system size and complexity increase; computationally intensive | Distributes training and execution across agents, enabling better handling of large-scale problems |
| Decentralization | Often rely on centralized and global system knowledge | Supports decentralized decision-making and coordination among agents |
| Partial Observability | Require full or simplified system visibility; struggle with uncertainty | Designed to operate effectively with incomplete or local observations using Dec-POMDP-based frameworks |
| System Heterogeneity | Limited to homogeneous models; struggle to accommodate diverse agents or tasks | Learns specialized policies per agent, handling heterogeneous roles and capabilities efficiently |

The classical approaches often face several limitations to solve resource allocation as summarized in the Table 2. Firstly, classical methods typically model and assume static or predictable environments, making them less adaptable to real-time changes, whereas MARL agents continuously learn and adjust to dynamic conditions. Secondly, scalability becomes a major challenge for traditional techniques as the size and complexity of the system grow, while MARL distributes learning across agents, enabling efficient performance even in large-scale settings. Additionally, classical methods generally rely on centralized control and full system visibility, which limits their applicability in decentralized and partially observable environments. MARL, in contrast, supports decentralized decision-making and performs well under local or incomplete information. Finally, conventional models often struggle to handle heterogeneous systems with varying agent roles or capabilities, while MARL can develop specialized policies tailored to diverse agent functions, allowing for more flexible and robust resource allocation.

These limitations arise some challenges underscore the need to move beyond classical approaches toward more advanced, AI-driven RAO methods, with a particular focus on MARL. In the following sections, we first introduce the foundational concepts of RL and the key principles of MARL. We then survey MARL solutions in addressing RAO problems, structured as follows: Section 4.2.1 explores solutions for continuous and rapidly changing issues, Section 4.2.2 examines approaches to decentralization and partial observability, Section 4.2.3 addresses scalability challenges, and Section 4.2.4 reviews strategies for managing heterogeneity in resources and objectives.

# 3 MARL Foundations

In this section, we begin by reviewing the essential concepts of RL, which form the basis for MARL. By understanding RL's approach to adaptive decision-making and

sequential optimization, we can see how these principles extend naturally to multi-agent environments. MARL leverages the decentralized nature of multiple interacting agents to overcome issues of scalability and enable real-time adaptability in resource allocation.

## 3.1 Reinforcement Learning as Optimization

RL is fundamentally suited to address the decision-making needs in RAO by enabling agents to learn optimal resource allocation policies through trial and error in dynamic and uncertain environments. Through interactions with the environment, agents can identify actions that maximize long-term rewards while adapting to immediate changes in resource demands or availability. This interactive learning process aligns well with RAO's goal of balancing long-term efficiency with real-time demands. The agent's Q-function estimates the expected cumulative reward for each action in a given state, guiding optimal actions that balance immediate and future outcomes. In RAO, this allows the agent to dynamically allocate resources based on the current system state and its expected impact on future performance Mao et al. (2016).

One of RL's key advantages in RAO is its capacity to adapt to changing conditions Vengerov (2007). Real-world RAO applications often involve unpredictable fluctuations, such as varying workloads in cloud computing or shifts in energy demand in smart grids. RL's continuous learning framework allows agents to adjust allocation strategies in real time, accommodating shifts in resource availability or demand. Importantly, RL also supports decision-making under uncertainty, allowing agents to explore allocation strategies even without full knowledge of future requirements or the actions of other agents.

Reinforcement Learning (RL) is a machine learning paradigm where an agent learns to make sequential decisions by interacting with an environment to maximize cumulative rewards over time Sutton and Barto (2018). In RL, the agent's objective is to learn a policy (a mapping from states to actions) that optimizes long-term rewards, adjusting its actions based on feedback from the environment.

### 3.1.1 Markov Decision Process

The RL framework is formally represented as a Markov Decision Process (MDP), which models the environment, agent actions, state transitions, and rewards Sutton (1988). This structure allows agents to dynamically adapt their actions based on observed states and obtained rewards, offering potential for real-time, flexible resource allocation in static or semi-static conditions. An MDP is defined by a tuple $(S, A, p, r, \gamma)$. Here, $S$ represents the set of possible states, while $A$ denotes the set of actions available to the agent. The transition dynamics are modeled by the probability function $p(s'|s, a)$, which describes the probability of moving from state $s$ to state $s'$ after taking action $a$. The immediate reward, $r(s, a, s')$, captures the reward received from transitioning from $s$ to $s'$ by action $a$. The discount factor $\gamma \in [0, 1]$ determines the relative importance of future rewards in the optimization process, weighting immediate rewards more heavily when closer to zero.

### 3.1.2 Learning Objective and Bellman Equation

In an MDP, the objective is to find an optimal policy $\pi^*$, that maximizes the expected cumulative reward, or return, over time. The return at time step $t$, denoted $R_t$, is expressed as the sum of discounted rewards $R_t = \sum_{i=0}^{\infty} \gamma^i r_{t+i+1}$. To guide its decisions, the agent relies on value functions, specifically the state value function $V_\pi(s)$ and the state-action value function $Q_\pi(s, a)$. The state value function, $V_\pi(s) = \mathbb{E}_\pi [R_t | S_t = s]$, represents the expected return when beginning in state $s$ and following policy $\pi$. The state-action value function, $Q_\pi(s, a) = \mathbb{E}_\pi [R_t | S_t = s, A_t = a]$, provides the expected return when starting from state $s$, taking action $a$, and then following policy $\pi$. These value functions are recursively defined by the Bellman equations, which relate the current value of a state to the expected value of future states and rewards. The Bellman equation for the state value function $V_\pi(s)$ is given by:

$$V_\pi(s) = \sum_a \pi(a|s) \sum_{s'} p(s'|s, a) \left[r(s, a, s') + \gamma V_\pi(s')\right]. \tag{2}$$

The state value function $V_\pi(s)$ can alternatively be expressed in terms of the action-value function $Q_\pi(s, a)$ as:

$$Q_\pi(s, a) = \sum_{s'} p(s'|s, a) \left[r(s, a, s') + \gamma \sum_{a'} \pi(a'|s')Q_\pi(s', a')\right]. \tag{3}$$

The policy $\pi$ aims to maximize this equation, yielding the optimal state-action value $Q^*(s, a)$ across all states and actions.

### 3.1.3 Deep RL

As environments become increasingly complex and involve more factors, traditional methods like SARSA and Q-learning struggle with scalability and stability. This has led to the development of advanced algorithms such as Deep Q Networks (DQN), which approximate the Q-value function with a neural network $Q(s, a; \theta)$, where $\theta$ represents the network parameters. To stabilize training, DQN introduces Experience Replay and a target Q-network to decouple updates. The target Q-value and the corresponding loss function are defined together as:

$$L(\theta) = \mathbb{E} \left[\left(r + \gamma \max_{a'} Q(s', a'; \theta') - Q(s, a; \theta)\right)^2\right]. \tag{4}$$

Enhancements like Double DQN Van Hasselt et al. (2016) and Dueling DQN Wang et al. (2016) improve performance by reducing overestimation bias and separating state values from action advantages.

Policy Gradient (PG) methods optimize the policy by directly maximizing the expected return $J(\theta)$ using Monte Carlo estimates. In the vanilla PG approach, the

gradient is given by:

$$\nabla_\theta J(\theta) = \mathbb{E}_\pi \left[ \sum_{t=0}^{T} \nabla_\theta \log \pi(a_t|s_t; \theta) G_t \right], \tag{5}$$

where $G_t$ is the cumulative reward from time step $t$ onward. Actor-Critic methods Konda and Tsitsiklis (1999) combine Policy Gradient (PG) methods with value function estimation, where the actor updates the policy parameters, and the critic evaluates the value function. Advantage Actor-Critic (A2C) Mnih et al. (2016) further reduces variance by using the advantage function, $A(s,a) = Q(s,a) - V(s)$, instead of the value function.

To maintain stability, Trust Region Policy Optimization (TRPO) Schulman (2015) and Proximal Policy Optimization (PPO) Schulman et al. (2017) constrain updates to prevent large deviations from the current policy. TRPO achieves this by limiting KL divergence while PPO uses a clipped objective.

## 3.2 Multi-Agent Reinforcement Learning

MARL extends the reinforcement learning framework to environments where multiple agents interact and learn concurrently within a shared space Bu et al. (2008). While standard RL methods focus on optimizing the cumulative reward of a single agent, MARL addresses the complexities arising from multiple autonomous agents learning and adapting simultaneously. This is especially relevant for RAO in large-scale, decentralized, and dynamic systems where agents must cooperate or compete to manage limited resources effectively.

### 3.2.1 Why MARL Match to Modern RAO

The capabilities of MARL make it highly effective for addressing the core challenges in modern RAO. Its decentralized training and execution paradigms enable agents to learn and act autonomously, reducing the computational bottlenecks and data flow issues inherent in centralized systems Ma et al. (2024). This scalability is critical for RAO applications that span large-scale networks, such as telecommunications or smart grids, where centralized solutions fall short.

Additionally, MARL's continuous learning processes allow agents to update their policies based on real-time feedback, making the approach highly adaptable to rapidly changing environments, such as fluctuating energy demands in power grids or dynamic workloads in cloud computing. In RAO settings that involve multi-agent systems, such as multi-robot teams or IoT networks, MARL also supports decentralized operations Zhou et al. (2023). This is especially useful under the CTDE paradigm, where agents rely on local observations and make decentralized decisions, addressing the need for distributed control in large-scale systems with limited information Charbonnier et al. (2022).

MARL is also well-suited to handle heterogeneous resources and diverse objectives. It allows for multi-objective optimization, enabling agents to balance local goals with

overarching system objectives—a crucial feature for systems requiring distinct allocation strategies for varied resources like bandwidth, energy, and memory Xiao et al. (2023). Moreover, MARL's support for stochastic games and Dec-POMDPs equips agents to make informed decisions in uncertain environmentsNguyen et al. (2020), a necessity in RAO tasks with unpredictable variables such as fluctuating network loads in telecommunications or variable energy supplies in smart grids.

By providing a structured and adaptable framework, MARL offers a decentralized, scalable, and responsive approach to the demands of dynamic and complex resource allocation environments Ning and Xie (2024). The following sections will delve into specific MARL algorithms and methods tailored to these challenges, examining their potential to meet the unique requirements of modern RAO applications.

### 3.2.2 Stochastic Games

In multi-agent systems (MAS), agents operate within cooperative, competitive, or mixed environments. These interactions are captured through the formalism of *Stochastic Games*, a generalization of Markov Decision Processes (MDPs) that accommodates multi-agent dynamics Hu and Wellman (2003). A stochastic game, or multi-agent MDP (MA-MDP), involving $N$ agents is defined by the tuple: $\langle S, \{A_i\}_{i=1}^N, T, \{r_i\}_{i=1}^N \rangle$, where $S$ represents the set of states, $A_i$ is the set of actions available to agent $i$, forming the joint action set $A = A_1 \times \cdots \times A_N$, $T : S \times A \times S \to [0,1]$ is the state transition function, and $r_i : S \times A \to \mathbb{R}$ is the reward function for each agent $i$.

In stochastic games, the state transitions depend on the joint action $\mathbf{a} = (a_1, \ldots, a_N)$ of all agents, which is critical for addressing dynamic RAO challenges. In fully cooperative settings, all agents share a single reward function, aligning with the goal of optimizing resource allocation for a common objective. In contrast, competitive settings introduce conflicting objectives (e.g., in zero-sum games where $r_1 + r_2 = 0$ for two agents), while mixed settings involve a combination of cooperative and competitive elements, providing flexibility to model various RAO scenarios.

### 3.2.3 Partially Observable MAS and Dec-POMDP

In practical multi-agent systems, agents frequently face *partial observability*, where each agent only has limited information about the environment's state, introducing additional complexity due to uncertainty. To address this, *Partially Observable Markov Decision Processes (POMDPs)* extend the MDP framework, allowing agents to make decisions based on partial and uncertain observations.

The multi-agent version, known as the *Decentralized Partially Observable Markov Decision Process (Dec-POMDP)*, is well-suited for decentralized RAO scenarios where agents operate based on localized information Oliehoek et al. (2016). A Dec-POMDP is defined by the tuple: $\langle S, \{A_i\}_{i=1}^N, T, r, \{O_i\}_{i=1}^N, O, N, \gamma \rangle$, where $S$ represents the set of environment states, and each agent $i$ has an action space $A_i$, forming the joint action space $A = A_1 \times \cdots \times A_N$ for $N$ agents. The state transition function $T : S \times A \times S \to [0,1]$ describes the probability of transitioning from state $s$ to $s'$ given the joint action $\mathbf{a} = (a_1, \ldots, a_N)$. The global reward function $r : S \times A \to \mathbb{R}$ provides feedback based on the joint actions. Each agent $i$ has an observation space $O_i$, and the joint observation

(a) Centralized Training and Centralized Execution (CTCE)

(b) Decentralized Training and Decentralized Execution (DTDE) (Wen et al., 2021)

(c) Centralized Training and Decentralized Execution (CTDE)
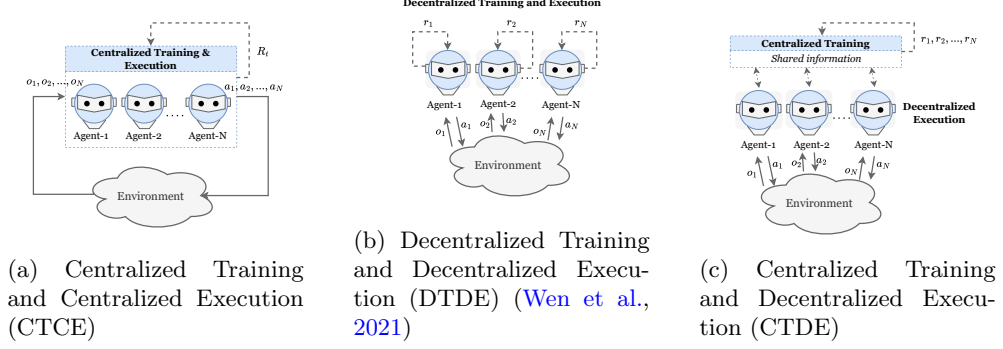
**Fig. 4**: Comparison of Training and Execution Paradigms in MARL: (a) CTCE is a fully centralized framework which combines all observations and actions of the agents into a joint observation-action space. (b) DTDE is a term for fully decentralized setting that treats all agents independently with their own observation, action, and reward. (c) CTDE framework has training in a centralized manner with information from other agent, then deploy the trained policy to each agent independently.

space is $O = O_1 \times \cdots \times O_N$. The observation function $O : S \times A \times O \to [0, 1]$ gives the probability of agent $i$ receiving observation $o_i$ given state $s$ and joint action $\mathbf{a}$. The number of agents is $N$, and $\gamma \in [0, 1]$ is the discount factor that determines the importance of future rewards.

Dec-POMDPs inherently support cooperative settings, aligning with RAO tasks that rely on collaboration for optimal resource distribution. Cooperative MARL enables agents to work collectively to achieve system-wide objectives, addressing modern challenges like scalability by distributing decision-making across agents. The decentralization inherent in Dec-POMDPs addresses the impracticality of central control in large-scale systems, while partial observability reflects the realistic limitations on agents' knowledge in decentralized RAO settings.

### 3.2.4 Learning Paradigms

The training of agents in MARL follows three main paradigms. Each paradigm presents unique benefits and trade-offs for modern RAO, depending on the scale, distribution, and interdependence of resources. Figure 4 illustrates the structural differences among these paradigms.

#### *Centralized Training and Centralized Execution (CTCE)*

In CTCE paradigm, a centralized controller determines the actions for all agents based on the global state, optimizing coordinated interactions. The controller utilizes a global policy $\pi_{\text{global}}(s)$ or a value function $Q_{\text{global}}(s, \mathbf{a})$ to select the optimal joint action $\mathbf{a}$, either by setting $\mathbf{a} = \pi_{\text{global}}(s)$ or by choosing $\mathbf{a} = \arg\max_{\mathbf{a}'} Q_{\text{global}}(s, \mathbf{a}')$. While CTCE ensures coordinated actions, scalability can be challenging as the number of agents increases, due to the exponential growth of the joint action space and the

communication demands required for real-time access to the global state. CTCE is best suited to controlled environments with relatively few agents, where centralized control remains practical, while it encounters limitations in larger or more decentralized RAO settings, such as distributed energy networks or cloud systems, where scalability and latency are significant considerations.

### *Decentralized Training and Decentralized Execution (DTDE)*

In DTDE, each agent operates independently, selecting actions based solely on its local observations. The policy $\pi_i(o_i)$ maps each agent's local observation $o_i$ to an action $a_i$ as $a_i = \pi_i(o_i)$, or alternatively, using a local value function $a_i = \arg\max_{a_i'} Q_i(o_i, a_i')$. This decentralized framework offers advantages in scalability and flexibility, enabling agents to operate autonomously. As such, DTDE is suitable for highly distributed settings, such as sensor networks or fully decentralized RAO scenarios. However, DTDE does not inherently provide coordinated actions, which can result in non-stationary environments where agents must adapt to the evolving policies of others. Independent learning (IL) methods, such as Independent Q-Learning (IQL) and Independent Proximal Policy Optimization (IPPO), are effective for large-scale, distributed tasks but may encounter challenges when coordination or shared objectives are essential, as seen in cooperative RAO contexts.

### *Centralized Training and Decentralized Execution (CTDE)*

CTDE combines the benefits of centralized training and decentralized execution. During training, agents have access to global state information, which allows for a centralized learning process. A global Q function $Q_{\text{global}}(s, \mathbf{a})$ or value function $V_{\text{global}}(s)$ is used to optimize the joint behavior of agents During execution, however, each agent follows its individual Q function $Q_i(o_i, a_i)$ or learned policy $\pi_i(o_i)$ based on local observations $o_i$, allowing for decentralized decision-making without access to the global state or knowledge of other agents' actions. Key approaches under CTDE include centralized critic and credit assignment methods. The centralized critic approach (e.g., MADDPG Lowe et al. (2017) and MAPPO Yu et al. (2022)) uses a centralized function during training to evaluate joint actions, while credit assignment methods like Value Decomposition Networks (VDN) Sunehag et al. (2018) and QMIX Rashid et al. (2020) decompose joint rewards to guide individual agents. CTDE's flexibility makes it ideal for complex RAO applications, such as managing resources in distributed cloud environments, where agents must independently make allocation decisions while aligning with overall system objectives.

## 4 RAO Leveraging MARL

MARL offers a powerful approach to tackling modern challenges in RAO, particularly where traditional methods fall short in scalability, dynamic adaptation, and decentralized decision-making. MARL extends reinforcement learning (RL) to involve multiple agents interacting within a shared environment, enabling a collaborative, decentralized approach to resource allocation that adapts to complex and changing conditions Bu et al. (2008). In this section, a summarization of the available MARL algorithm for

different application is discussed, then followed by the primary challenges addressed in that particular field and application.

## 4.1 MARL for RAO in Different Application Fields

This section explores how MARL serves as a powerful framework for addressing RAO problems across various domains. We begin by discussing how Reinforcement Learning (RL) functions as an optimization technique capable of handling dynamic and uncertain environments. Building on this foundation, we highlight the advantages of MARL in solving complex RAO tasks that involve multiple decision-making agents operating in decentralized or partially observable settings. We then present how MARL has been applied to RAO challenges in diverse real-world fields such as Telecommunication Network and IoT, Energy, Distributed Computing, Transportation, and Manufacturing. Finally, we identify and analyze key challenges commonly encountered in applying MARL to RAO, including adaptability, coordination, scalability, and heterogeneity in resources (see summary in Table 3).

### 4.1.1 Telecommunications and Network Management

#### *Vehicular Network*

Resource allocation in vehicular networks is a critical challenge due to the highly dynamic and decentralized nature of vehicular environments. With the growing demand for intelligent transportation systems, efficient allocation of communication, computing, and storage resources is essential to ensure low latency, high reliability, and optimal system performance. The integration of Multi-Agent Reinforcement Learning (MARL) has emerged as a promising solution for enabling distributed decision-making and adaptability in such dynamic settings. MARL allows vehicles and network nodes to learn cooperative strategies in real-time, handle partial observability, and respond effectively to time-varying network conditions. The reviewed papers demonstrate how MARL significantly enhances adaptability in vehicular network environments. Mostly, including Ji et al. (2023); Li et al. (2022); Parvini et al. (2023) incorporate adaptive MARL algorithms such as MAD3QN and MADDPG to dynamically respond to changing network conditions and user demands. Studies like Seid et al. (2021); Zhang et al. (2020); Wu et al. (2020) further apply these models to UAV-enabled IoT, secure communications, and traffic light control, highlighting the algorithms' ability to handle spatio-temporal variability. Distributed and cooperative setups in Hu et al. (2020) and Zhou et al. (2023) exemplify the importance of decentralized decision-making in real-time resource management. Although Zhao et al. (2023) and Yin and Yu (2021) do not explicitly state adaptability, they adopt advanced MARL frameworks capable of operating in heterogeneous and dynamic vehicular networks. Overall, these works underscore the strength of MARL in providing flexible, context-aware solutions for resource allocation, outperforming static or rule-based conventional methods.

#### *Mobile Network*

The reviewed studies in mobile network environments, particularly under the complexity of mobile network management and heterogeneous architectures. Du et al. (2022)

**Table 3**: MARL for RAO in different application and its primary challenge: A = Adaptability, PO = Partial Observability, LS = Large Scale, H = Heterogeneity

| Field | Application | MARL Algorithm | A | PO | LS | H |
|---|---|---|:---:|:---:|:---:|:---:|
| Telecommunication Network and IoT (Section 4.1.1) | Vehicular Network | MADQRL MADDPG MAD3QN | ✓ | ✓ | ✓ | - |
| | Mobile Network | Graph Based MADDPG MAPPO | ✓ | - | - | ✓ |
| | Wireless Network | MADQRL MAPPO | - | ✓ | ✓ | ✓ |
| | Computer Network | MADQRL | ✓ | - | ✓ | - |
| | IoT Network | MAPPO | - | ✓ | ✓ | - |
| Energy (Section 4.1.2) | Microgrid | MATRPO MADDPG MATD3 | ✓ | - | - | ✓ |
| | Smart grid | MAPPO | ✓ | - | - | |
| | Renewable Energy | MAD3QN MATD3 | ✓ | - | - | - |
| | Power Distribution | Graph Based MADDPG MATD3 | ✓ | ✓ | ✓ | - |
| | Home Energy Managemet | MADQRL | ✓ | - | ✓ | - |
| Distributed Computing (Section 4.1.3) | Satellite Edge Computing | MAAC | - | - | ✓ | - |
| | Mobile Edge Computing | MAD3QN MADDPG MAAC MAPPO | ✓ | - | ✓ | ✓ |
| | Fog Computing | MAAC | ✓ | - | - | - |
| | Vehicle Edge Computing | MAD3QN MAAC MADDPG MAPPO | ✓ | ✓ | - | ✓ |
| | Vehicle Fog Computing | MAAC | ✓ | ✓ | - | ✓ |
| Transportation (Section 4.1.4) | Traffic Management | MADDPG MAAC MADQRL Graph Based | ✓ | ✓ | ✓ | - |
| | Autonomous Vehicles | MADDPG MATD3 MAAC MADQRL Graph Based | ✓ | ✓ | ✓ | - |
| | Vehicle Allocation | MADDPG | - | - | ✓ | - |
| Manufacturing (Section 4.1.5) | Flexible Job shop | MAPPO Graph Based | ✓ | - | - | - |
| | Cognitive Manufacturing | Graph Based | ✓ | - | - | - |

proposed a graph-based MARL (GA-Net MARL) framework for dynamic resource management in 6G in-X subnetworks, emphasizing the capability of graph structures to capture inter-agent relationships and support scalability. Allahham et al. (2022) addresses the challenges of network selection and resource allocation in multi-RAT (Radio Access Technology) networks using MADDPG, enabling collaborative decision-making in diverse environments. Meanwhile, Kim and Lim (2021) leverages MAPPO for end-to-end network slicing, facilitating efficient and adaptive management of resources across network segments. Collectively, these works underline MARL's ability to handle the heterogeneity, real-time demands, and large-scale coordination challenges present in next-generation mobile networks.

### Wireless Network

The integration of MARL into wireless networks demonstrates promising capabilities for optimizing dynamic and distributed operations. Naderializadeh et al. (2021) employs MADQRL to manage wireless resources in a decentralized fashion, effectively tackling the coordination challenges among multiple agents in large-scale environments. In a complementary approach, Guo et al. (2020) proposes a MAPPO-based framework to jointly optimize handover control and power allocation, highlighting MARL's strength in handling coupled decision-making problems in wireless networks. Together, these studies showcase MARL's potential to outperform classical approaches by enabling adaptive, scalable, and cooperative solutions for complex wireless resource allocation scenarios.

### Computer Network

In computer network application, You et al. (2020) leverages a fully distributed MADQRL framework for packet routing, enabling each agent to make routing decisions independently while adapting to network dynamics. This design reduces reliance on centralized control and enhances scalability. Similarly, Suzuki et al. (2022) adopts a cooperative MADQRL approach for dynamic virtual network allocation under fluctuating traffic demands, demonstrating how MARL can effectively respond to variability in network load while maintaining efficient resource use. Both studies underscore MARL's ability to enable decentralized, adaptive, and resilient solutions for complex computer networking problems.

### IoT Network

In the IoT Network, Xiao et al. (2023) proposes a MAPPO-based multi-agent deep reinforcement learning framework to address resource allocation challenges in large-scale IoT networks, specifically tailored for ultra-reliable low-latency communication (URLLC) scenarios. Their approach enables decentralized agents to make real-time decisions under strict QoS requirements while coordinating effectively to handle the scale and complexity of controllable IoT systems. This work highlights how MARL, particularly with MAPPO, can meet the stringent demands of emerging IoT applications that require both scalability and reliability in dynamic environments.

### 4.1.2 Energy

*Microgrid*

Recent research has increasingly leveraged MARL for efficient and scalable energy management in microgrid systems. Xu et al. (2024) proposes a hierarchical trust-region MARL framework (MATRPO) to optimize operations across interconnected multi-energy microgrids, focusing on collaboration and trust-aware learning. Abid et al. (2024) develops a novel multi-objective optimization strategy employing MADDPG to enhance planning decisions for microgrid resource allocation, balancing multiple operational goals. Zhang et al. (2023) introduced a distributed control architecture using MATD3 for real-time energy management, enabling flexible coordination among microgrids with varying energy types. Complementarily, Jendoubi and Bouffard (2023) designs a hierarchical MARL model using MADDPG to support layered decision-making, improving learning efficiency and coordination in complex microgrid systems. These works collectively illustrate the diverse applications of MARL techniques in enhancing energy distribution, planning, and real-time control within smart microgrid infrastructures.

*Smart Grid*

In smart grid applications to enhance the efficiency and autonomy of energy systems. Kumari et al. (2024) introduces a decentralized residential energy management system leveraging Deep Q-Network-based MARL (MADQN), focusing on distributed decision-making for home energy optimization. Their approach enables agents to learn optimal consumption strategies while ensuring grid stability. Meanwhile, Roesch et al. (2020) applies MAPPO in an industrial smart grid context, where multiple agents coordinate to manage energy flows dynamically and adaptively, reflecting the increasing complexity of industrial power systems. These contributions underscore the effectiveness of MARL in addressing decentralized, real-time control challenges in smart grids.

*Renewable Energy*

Recent advances in MARL have significantly contributed to optimizing renewable energy integration across diverse systems. Jayanetti et al. (2024) proposes a MARL framework based on Multi-Agent Actor-Critic (MAAC) for renewable energy-aware workflow scheduling across distributed cloud data centers, aiming to improve energy efficiency and computational performance. Shen et al. (2022) develops a MAD3QN-based optimization framework for building energy systems, incorporating renewable sources to balance comfort and energy costs. Chen et al. Chen et al. (2022) introduces a physics-shielded MARL method utilizing MATD3 for active voltage control in photovoltaic and battery-integrated grids, enhancing both safety and operational stability. These studies showcase the versatility of MARL in addressing the complexities of renewable energy systems.

*Power Distribution*

In power distribution systems, MARL has emerged as a powerful tool for decentralized voltage and reactive power control. Hu et al. (2024) proposes a graph-based

MARL approach using a decentralized training and decentralized execution (DTDE) framework for Volt-VAR control, demonstrating scalability and coordination among distributed agents. Wang et al. Wang et al. (2021) applies both MADDPG and MATD3 algorithms to enhance active voltage regulation performance in distribution networks, improving adaptability in dynamic environments. Sun and Qiu (2021) introduces a two-stage Volt/VAR control method using MADDPG for active distribution networks, combining global planning and local reactive power adjustment. These works illustrate the effectiveness of MARL in enabling autonomous, robust, and cooperative voltage control strategies in modern power grids.

### Home Energy Management

In the application of Home Energy Management, MARL has been widely explored to optimize distributed control and enhance residential energy flexibility. Charbonnier et al. (2022) proposes a scalable MARL framework based on MAQRL to manage distributed energy resources while preserving user comfort and ensuring scalability across numerous households. Xu et al. Xu et al. (2020) introduces a data-driven home energy management method employing MADQRL, which efficiently adapts to dynamic consumption patterns and uncertain energy generation in smart homes. Similarly, Ahrarinouri et al. (2020) utilizes MADQRL to enable collaborative energy management among residential buildings, achieving significant improvements in both cost reduction and peak load management. These studies demonstrate the capability of MARL approaches to deliver intelligent, decentralized, and adaptive energy management solutions in residential environments.

## 4.1.3 Distributed Computing

### Satellite Edge Computing

In the domain of distributed computing, especially satellite mobile edge computing (SMEC), managing large-scale task offloading across multiple dynamic nodes is a complex challenge. Zhang et al. (2024) addresses this by proposing a Multi-Agent Actor-Critic (MAAC) based collaborative optimization framework. The approach is designed to be scalable and adaptable, enabling decentralized satellite nodes to efficiently learn offloading strategies under highly dynamic and resource-constrained conditions. By coordinating multiple agents in a shared environment, the solution effectively tackles the scalability issues typical in SMEC scenarios, making it promising for future large-scale satellite-enabled computing systems.

### Mobile Edge Computing

The MARL has emerged as a powerful approach for addressing complex challenges in Mobile Edge Computing (MEC), particularly in large-scale environments where resource allocation, task offloading, and adaptability to dynamic conditions are critical. The following papers highlight the challenges of large-scale resource allocation and adaptability in MEC. Liu et al. (2024) focuses on computation rate maximization for SCMA-aided edge computing in IoT networks with the MAD3QN algorithm, addressing scalability and adaptation in dynamic environments. Gao et al. (2023) and Gao

et al. (2022b) tackle large-scale cooperative task offloading and resource allocation in heterogeneous MEC systems, utilizing the MAAC algorithm to manage scalability and adaptability in diverse conditions. Gao et al. Gao et al. (2023) proposes Com-DDPG for task offloading in MEC systems for the internet of vehicles, emphasizing the ability to adapt to information-communication-enhanced environments. Zhao et al. (2022) presents MATD3 for task offloading in UAV-assisted MEC, focusing on scalable and adaptable solutions for mobile edge networks. Cao et al. (2020) applies MADDPG to address multichannel access and task offloading challenges in Industry 4.0, emphasizing large-scale adaptability. Lastly, Wu et al. (2023) uses MAPPO for minimizing completion delay and energy consumption in MEC-based Industrial IoT (IIoT) systems, emphasizing scalability in handling large-scale IIoT environments.

### Fog Computing

In the application of Fog Computing, Jain and Kumar (2023) focuses on QoS-aware task offloading in a fog computing environment, utilizing the MAAC algorithm to address resource allocation and task management challenges in distributed, resource-constrained fog networks.

### Vehicle Edge Computing

Kang et al. (2023) presents a MAPPO-based approach for cooperative UAV resource allocation and task offloading in hierarchical aerial computing systems. Ju et al. (2023) proposes a joint secure offloading and resource allocation strategy for vehicular edge computing networks using MADDQN, addressing both security and resource management challenges. Zhu et al. (2020) utilizes the MAAC algorithm for vehicular computation offloading in IoT environments, emphasizing task allocation in vehicular edge computing systems. Zhang et al. (2021) combines adaptive digital twin technology with MADDPG to optimize resource management in vehicular edge computing and networks. These studies highlight the importance of scalability, adaptability, and security in resource allocation and task offloading for vehicle edge computing systems.

### Vehicle Fog Computing

Wei et al. (2023) explores many-to-many task offloading in vehicular fog computing using the MAAC algorithm, addressing the challenge of efficiently managing resource allocation and task offloading between multiple vehicles and fog nodes. Gao et al. Gao et al. (2022a) proposes a fast adaptive task offloading and resource allocation approach in heterogeneous vehicular fog computing systems, also utilizing the MAAC algorithm to enhance the efficiency and adaptability of resource management in dynamic vehicular environments. Both studies emphasize the scalability and flexibility required for efficient resource allocation in vehicular fog computing.

### 4.1.4 Transportation

#### Traffic Management

The following research works explore the application of MARL in traffic management. Zhang et al. (2024) presents MARLens, a visual analytics approach to understanding

traffic signal control using MADDPG.Chen et al. (2023) applies MAAC for highway on-ramp merging in mixed traffic to enhance traffic flow. Zeynivand et al. (2022) implements MADQRL for traffic flow control, optimizing vehicular movement across networks. Wang et al. (2020) introduces STMARL, a spatio-temporal MARL approach for cooperative traffic light control using a graph-based method. Wang et al. (2020) address large-scale traffic signal control with MADQRL, optimizing traffic management in urban settings. Wu et al. (2020) utilizes MADDPG for urban traffic light control in vehicular networks, focusing on multi-agent coordination to improve traffic efficiency. These studies demonstrate the versatility and effectiveness of MARL algorithms in optimizing various aspects of traffic management.

### Autonomous Vehicle

MARL is applied to autonomous vehicles for optimizing decision-making, coordination, and resource allocation in various driving environments. MARL approaches enable autonomous vehicles to interact, cooperate, and make decisions in complex, dynamic settings, improving safety, efficiency, and traffic flow. Antonio and Maria-Dolores Antonio and Maria-Dolores (2022) applies MATD3 to manage connected autonomous vehicles at intersections, optimizing vehicle coordination and traffic flow. Jiandong et al. (2021) explores UAV cooperative air combat maneuvers, using MAAC for decision-making in autonomous vehicle control scenarios. Chen et al. (2021) leverages graph neural networks combined with reinforcement learning for multi-agent cooperative control of connected autonomous vehicles, optimizing vehicle coordination in complex networked environments. These studies demonstrate the potential of MARL for enhancing the coordination and decision-making capabilities of autonomous vehicles in both traffic and cooperative tasks.

### Vehicle Allocation and Routing

MARL has shown promising results in the application of autonomous vehicles, particularly for resource allocation and optimization in complex transportation systems. By enabling vehicles to interact and cooperate with each other, MARL approaches can improve routing, task allocation, and overall operational efficiency in autonomous vehicle systems. Ren et al. (2022) applies MADDPG to optimize vehicle routing in supply chain management, focusing on the allocation and efficient routing of vehicles. Their approach integrates route recorders to enhance coordination and decision-making, ensuring that vehicles are efficiently assigned tasks and follow optimal routes, thereby improving supply chain logistics and transportation efficiency. This study showcases the potential of MARL in optimizing vehicle allocation and routing, with a focus on practical applications in supply chain management.

### 4.1.5 Manufacturing

In manufacturing systems, MARL can significantly improve RAO task by enabling decentralized decision-making among autonomous agents, enhancing flexibility, efficiency, and adaptability in complex production environments.

*Flexible Job Shop*

MARL has emerged as an effective approach to solve RAO problems in complex manufacturing process, especially in flexible job shop scheduling. In such systems, multiple agents must make real-time, adaptive decisions to allocate resources, optimize production workflows, and meet dynamic constraints, which is essential for increasing efficiency and reducing costs. Jing et al. (2024) utilizes a graph-based MARL approach combined with Graph Convolutional Networks (GCN) to optimize flexible job shop scheduling, enhancing the adaptability of agents in handling complex scheduling tasks. Heik et al. (2024) proposes MAPPO to dynamically allocate manufacturing resources, addressing the need for adaptive decision-making under changing production conditions. Zhang et al. (2023) introduces DeepMAG, which integrates deep reinforcement learning with multi-agent graphs for flexible job shop scheduling, allowing agents to adapt to evolving task requirements and machine statuses. Liu et al. (2023) combines deep reinforcement learning and a multi-agent system to dynamically schedule re-entrant hybrid flow shops, accounting for worker fatigue and skill levels, thus improving adaptability in labor resource allocation. Finally, Zhang et al. (2022) focuses on dynamic job shop scheduling using MAPPO, which enables multi-agent manufacturing systems to adapt to real-time changes and uncertainties in the production process. These studies highlight the crucial role of adaptability in optimizing flexible job shop scheduling through MARL, especially dynamic manufacturing environments.

*Cognitive Manufacturing*

Cognitive manufacturing, which integrates AI and machine learning, has significant potential for improving decision-making, resource optimization, and process efficiency in manufacturing systems. MARL plays a key role in cognitive manufacturing by enabling multiple agents to work together in a decentralized manner to learn and adapt to various manufacturing tasks and environmental changes. Zheng et al. (2021) addresses this challenge by proposing an industrial knowledge graph-based MARL approach for cognitive manufacturing. Their approach enables agents to dynamically learn and adapt to different manufacturing tasks through the use of knowledge graphs, which helps in representing complex relationships between various system components.

## 4.2 MARL as a Solution for Modern RAO Challenges

The complexity of modern RAO has led to the extension of RL to multi-agent settings, enabling each agent to make localized and distributed resource allocation decisions. MARL offers a range of paradigms that effectively address the specific challenges in RAO. By extending traditional reinforcement learning to handle multiple agents within a shared environment, MARL introduces solutions that leverage centralized, decentralized, and hybrid learning paradigms. This flexibility enables MARL to meet scalability, adaptability, coordination, and privacy needs in RAO, making it a valuable framework for modern, complex systems. In this subsection, we discuss further four primary challenges and how it is solved

### 4.2.1 Adaptability in Continuously Changing Environments

In dynamic and uncertain environments, RAO must handle continuous changes in resource demands and availability that may often change unpredictably. An RL algorithm can be extended to multi-agent settings by simply combining its observations and actions from multiple agents. It constructs a single RL model with a larger amounts of inputs and outputs in a CTCE or fully-centralized framework. By using the CTCE framework, it allows the central controller to quickly adapt the overall strategy by recalculating optimal allocations based on current conditions. This agility is especially valuable in RAO contexts with fluctuating demands, enabling the system to maintain efficient allocations even as conditions evolve.

In Jain and Kumar (2023), a fully centralized MARL algorithms has been evaluated based on three different DRL methods such as DQN, DDPG, and SAC. It focused on handling the unpredictability of tasks and maintain the Quality of Service (QoS) requirements of users by considering a highly dynamic requirements, such as: end-to-end latency, energy consumption, task deadline, and priority. The algorithms are trained offline in a resource rich cloud data center. The SAC algorithm outperforms other baseline techniques in terms of time, energy consumption, utility, execution rate, and aging. A decentralized framework may struggle to adapt to these shifts due to limited information access, resulting in slower response times and potential under-utilization or overloading of resources.

In resource allocation, multiple agents may compete for limited resources, and optimal allocation requires that each agent's actions are highly coordinated to efficiently achieve the same goal. Centralized systems have access to all the agents' states, actions, and rewards, enabling the learning algorithm to optimize resource allocation globally rather than locally. This framework leads to more optimal resource allocation across the entire system since decisions are made with a comprehensive view of the environment. With the full access to all the agents' information and the entire system's state, a centralized system can make more well-informed decisions. This can be particularly beneficial in a cooperative resource allocation when global knowledge is necessary to avoid any over-use or under-use of resources. Decentralized approaches have a partially observable settings that make it struggles with this issue and leading to conflicts or suboptimal outcomes. A deep reinforcement learning algorithm, called Compounded-Action Actor-Critic (CA2C), has been evaluated to address the trajectory planning problem for the cellular Internet of UAVs in Hu et al. (2020). The CA2C algorithm is capable of effectively managing agents with complex actions, which involve both continuous and discrete variables. This work implements Deep reinforcement learning algorithms that is well-suited for determining optimal policies for agents in MDPs model that have high-dimensional state spaces. All states and actions are jointly processed in a single central controller. They evaluate the proposed algorithm by comparing it against four baseline algorithms. The proposed CA2C algorithm was shown to outperform four benchmark algorithms in terms of AoI minimization.

A fully centralized systems struggle with the scalability issue as the number of agents increases. The central controller must process all agents' states, actions, and rewards in a single computation resource leading to a dramatic increase in computational complexity and memory requirements. This can make centralized MARL

impractical for large-scale resource allocation problems. From the reliability point of view, a centralized system is vulnerable to failures in the central node controller. If the central node fails, the entire system can be disrupted, which reduces the system's fault tolerance and resilience. This can be particularly risky in mission-critical resource allocation tasks. Since all agents' information must be shared with the central controller, fully centralized systems pose privacy and security risks. In resource allocation scenarios involving sensitive data (e.g., personal data in healthcare Guindo et al. (2012), financial data in blockchain-based resource management Yánez et al. (2020)), these concerns may make centralized MARL less viable. Therefore, this settings are not widely used in RAO research field.

### 4.2.2 Coordinating Resources with Partial Observability

In many RAO application, each agent can only observe a limited view of the overall environment (e.g., local resource availability or neighboring agents' actions). For instance, in a network setting, an agent may only know its local bandwidth but not the network congestion affecting other agents. This fragmented view makes it challenging for agents to make well-informed decisions about resource allocation, often leading to suboptimal use of resources. Partial observability leads to scenarios where agents might allocate resources based on incomplete or outdated information, which can create conflicts, over-allocations, or wasted resources. For instance, an agent might allocate a resource already claimed by another, or under-utilized resources that could have been shared more effectively.

RAO settings often demand that agents make decisions with limited local information, particularly in IoT networks, multi-robot systems, and other distributed architectures. To optimize resource allocation, agents must coordinate their actions with limited information. However, partial observability inherently limits their ability to understand how their decisions impact others, especially in large systems where agents' actions can have ripple effects across the environment. Inadequate coordination can result in inefficient policies where agents inadvertently work against each other. Addressing partial observability in resource allocation requires a combination of these techniques to create a more cohesive MARL framework, enabling agents to coordinate effectively despite limited information. This way, agents can make more informed decisions, improving both local and global resource utilization. For example, in a grid computing system, multiple agents might allocate tasks to the same computing node, overloading it and slowing down task processing, while other nodes remain idle. This lack of coordination wastes computational resources and increases processing times. Here, **DTDE** is highly applicable, as it allows agents to make autonomous decisions based on local observations, thus aligning well with environments where agents have constrained or partial access to global states. The independence afforded by DTDE is crucial in applications where maintaining privacy and minimizing communication costs are essential, such as decentralized task allocation in manufacturing systems Zheng et al. (2021).

In an interconnected multi-energy microgrid optimization, Zhang et al. (2023) proposed solution to the decentralization problem in microgrid system. The authors implemented MADRL algorithm with an attention mechanism added to the centralized

critic to meet local customized energy demands in a form of decentralized execution. In this work, CTDE framework is used to maintain global optimization performance with the coordination of each agents. By using this framework, agents are trained with access to a centralized view of the environment, allowing them to learn optimal strategies that account for the whole system's needs. During execution, agents act based on their local observations, but they are guided by policies shaped by this global perspective.

However, in cases where some level of centralized oversight is feasible, CTCE can provide fully coordinated solutions. CTCE is ideal for small-scale RAO tasks with a limited number of agents requiring precise synchronization, such as coordinated robotic tasks in structured environments Jain and Kumar (2023). In this setting, CTCE enables the system to optimize resource use by leveraging a complete, centralized understanding of all agents' actions and states.

### 4.2.3 Dealing with Large-Scale Systems

In centralized or semi-centralized approaches, scaling up the number of agents increases the computational complexity and memory requirements due to the prodigious amount of state and action information that needs to be processed. The complexity grows exponentially as the number of agents increases, making it impractical for real-time resource allocation. Besides, In a DTDE framework, each agent is trained and operates independently, relying only on local observations without needing centralized information Jiandong et al. (2021); Wu et al. (2020). In other words, the agents are trained to maximize their local rewards and optimize their policies independently. This independence feature reduces computational overhead, as each agent handles its learning and decision-making based on local states, which is more scalable than a central agent tracking the states and actions of all others. With more agents, there is a risk that individual decision quality may degrade if agents lack sufficient information or coordination. In large-scale systems, maintaining high decision quality is critical to ensure efficient use of resources across all agents Gao et al. (2023). Through local interactions and learning policies tailored to specific environments, agents can still achieve near-optimal decisions at scale, even without full knowledge of other agents' actions. This framework can maintain decision quality by focusing on improving each agent's local policy, which contributes to a more stable system-wide resource allocation.

Scalability is a core challenge in large-scale resource allocation systems, where handling numerous agents and their interactions becomes increasingly complex. This problem emerge as the system size, number of agents, or complexity of the resource environment grows. Specifically, scalability issues arise when the algorithms or strategies used to allocate resources cannot efficiently handle an increase in agents or resource demands. In RAO, each agent must choose actions that contribute to optimal resource allocation. As the number of agents or resources increases, the combined action space grows exponentially, leading to a massive increase in the number of possible allocation combinations. Managing this vast space is computationally challenging, and traditional methods can struggle to identify optimal or even near-optimal solutions within a feasible time frame. This problem is particularly pronounced in real-time RAO scenarios, such as network bandwidth allocation or power distribution,

where decisions must be made quickly. The impact that can be occurred is the agents might redundantly allocate resources, over commit shared resources, or even leave some resources underutilized due to lack of information about other agents' decisions without effective. This results in inefficiencies and potential bottlenecks in resource distribution. Leveraging MARL frameworks, either DTDE or CTDE, can help agents to learn efficient policies in large-scale and partially observable environments without a central controller.

To solve scalability in RAO problem, DTDE scales well with a large number of agents. Each agent learns independently, allowing the system to handle more agents without requiring centralized coordination or training. This is particularly beneficial in resource allocation scenarios where there are many resources and agents in distributed energy systems or communication networks. Since agents train and execute independently, DTDE systems do not require the sharing of sensitive information among agents or a central controller. This feature can be crucial in situations like network resource allocation, where privacy and security concerns are high. DTDE allows agents to adapt to local environmental changes or constraints without waiting for a centralized update or communication. This can lead to more dynamic and responsive resource allocation, as agents can adjust in real-time based on localized resource availability and demand. In DTDE, agents are trained to act independently without a centralized controller, reducing the computational burden and making the system scalable. In Gao et al. (2023), a task offloading based on MARL has been implemented for Information-Communication-Enhanced Mobile Edge Computing for the Internet of Vehicles. This work evaluates LSTM network and a BRNN, Actor Critic MARL in a fully decentralized setup and compares all of the algorithms performance.

**Graph based MARL** has been evaluated to solve problem in power distribution using DTDE framework in Hu et al. (2024). The authors proposed algorithm divides the power distribution system into several regions, each region treated as an agent. Then an MARL was designed by employing Hierarchical Graph Recurrent Network (HGRN) structure that combines the advantages of Hierarchical Graph Attention (HGAT) and Gated Recurrent Unit (GRU) enables the communication between heterogeneous agents. Another graph based MARL has been implemented in Zheng et al. (2021) by using industrial knowledge graph (IKG)-based multi-agent reinforcement learning. It is designed to solve the robot task allocation and completion problem in a manufacturing network.

MARL's capacity for decentralized real-time adaptation is a key advantage here, particularly in **CTDE** framework, where agents are trained with global insights, yet operate with decentralized policies in real-time. For instance, in energy management systems, CTDE enables agents to adapt to fluctuating demands while maintaining overall grid stability Guo et al. (2020). By updating policies based on local observations and previously learned global strategies, agents can quickly react to dynamic conditions without centralized intervention. Furthermore, CTDE's use of a centralized critic during training mitigates the non-stationary issue that arises from agents learning simultaneously. This is particularly effective in RAO applications like edge computing, where agents (servers or virtual machines) optimize resource allocation across

fluctuating workloads Wu et al. (2023). By combining centralized training with decentralized adaptability, CTDE balances the need for coordination with the flexibility to handle rapid environmental changes.

By training with centralized information but executing with decentralized policies, the system can scale to a large number of agents and environments. During training, agents can learn to cooperate and avoid conflicts or inefficient resource usage. Decentralized execution allows agents to quickly adapt to local changes in demand or resource availability without waiting for global updates. Another benefit is that CTDE reduces communication overhead during execution, which is crucial in systems where bandwidth or computational power is constrained.

One type of Credit Assignment is developed using Value Decomposition (VD) method. The common algorithm that has been used in RAO problem is QMix. It has been evaluated to solve a problem for vehicular cloudlet in automotive industry Ahmed et al. (2023). A global critic (centralized function) is used to evaluate the actions of all agents collectively. The global reward is computed based on the combined agent's performance, and the centralized system updates a joint action-value function. Currently, VD has been rarely used in RAO compared to Centralized Critic since it has more abstract cooperative strategy represented by decomposing the joint action-value function $Q_{tot}$ into individual action-value functions $Q_i$ for each agent. The decomposition ensures that the sum or combination of the individual values approximates the global value function. However, there is no guarantee that it will converge to global optimal point during the training phase.

In cloud environments, CTDE can be applied to allocate computing resources (e.g., CPU, memory) among multiple virtual machines or applications. Agents (representing servers or VMs) learn to allocate resources based on system load and application requirements, optimizing for performance and cost. In Wu et al. (2023), a task offloading optimization problem for edge computing based Industrial IoT (EIIoT) infrastructure has been solved using MARL algorithm. An actor-critic (AC) DRL framework is adopted to construct a suitable lightweight offloading decision system and optimize the joint completion delay and energy consumption in the EIIoT.

CTDE MARL can be used to allocate spectrum resources dynamically between different users or devices in a 5G or IoT network. Agents (representing base stations or users) learn to allocate bandwidth efficiently based on energy consumption Meng et al. (2020); Nasir and Guo (2019), traffic demand and channel quality Kim and Lim (2021).

In energy management systems, CTDE MARL can help to allocate power from renewable sources among consumers or storage devices. Agents (representing homes, factories, or storage units) learn to request power optimally to balance supply and demand, minimize waste, and ensure grid stability. MARL has been evaluated an optimization of power allocation for home energy management using MAPPO algorithm Guo et al. (2020). In fact, MAPPO is also built upon centralized critic network and it shows competitive performance to solve the task in energy management.

### 4.2.4 Handling Heterogeneity in Resources and Objectives

In multi-agent systems, agents must balance diverse resource types, availability, and possibly conflicting objectives. Resources often vary widely in their characteristics, utility, and constraints, leading to heterogeneity in their availability, value, and compatibility with different tasks. For instance, in cloud computing, resources like CPU, memory, storage, and GPU have different capacities, costs, and efficiencies, and not all tasks require or benefit from each resource in the same way. In smart grids, energy sources might vary (e.g., solar, wind, fossil fuel), each with unique costs, emissions, and reliability levels. Matching the right resource to the right demand while maintaining system efficiency becomes complex in such diverse resource pools. Resource heterogeneity can create allocation inefficiencies and under-utilization. Allocating a resource poorly suited to a task can lead to degraded performance, over-utilization of high-value resources, and an imbalance where some resources are over-consumed while others are underutilized. Heterogeneity in resources and objectives adds complexity to the task of matching available resources with the diverse needs of agents. Poor matching leads to inefficiency, where high-demand or specialized resources are wasted on low-priority or incompatible tasks. This can reduce overall performance, degrade system stability, and prevent high-priority tasks from accessing the necessary resources.

Many RAO problems involve multiple types of resources, each with unique allocation goals. **CTDE** is well-suited for these heterogeneous environments, as it allows centralized training to capture complex interdependencies between diverse resource types. For example, in cloud computing, agents trained with CTDE learn to balance between CPU, memory, and bandwidth resources Wu et al. (2023), optimizing across multiple metrics such as energy consumption, latency, and throughput.

In situations with diverse agent goals or conflicting objectives, **Credit Assignment** within CTDE is particularly valuable. Credit assignment decomposes the joint action-value function into individual agent rewards, enabling each agent to receive feedback on its contribution to the overall objective. Value decomposition methods, such as QMix, can facilitate nuanced resource allocation in scenarios like vehicular cloudlets, where agents (vehicles) must manage bandwidth and energy resources to optimize communication Ahmed et al. (2023).

A centralized critic (by using actor-critic methods) or a centralized policy network is used to optimize the overall objective, accounting for interactions between agents as in Wei et al. (2023). During training, a centralized critic has access to the global state of the system (e.g., overall network load, resource availability, and actions of all agents). This allows the agents to learn how their individual actions affect overall system performance. Agents can share experiences and learn from each other, which speeds up training and improves coordination. Centralized critics approaches address the non-stationary problem caused by the changing behaviors of agents. These help to stabilize learning by using a shared critic function to evaluate the actions of multiple agents, reducing the complexity of dealing with multiple evolving policies. Once the training is complete, agents operate independently based only on their local observations. They no longer have access to the full state of the environment or the actions of other agents. This decentralized execution is scalable to large environments with multiple agents, as each agent can act autonomously without relying on global information in real time.

# 5 Available RAO Simulators for Benchmarking

In MARL, a simulator or environment for benchmarking is developed based on Open-AI gymnasium (gym) library as their main framework to be integrated with any common reinforcement learning algorithms. Here, we listed some publicly available benchmarks for RAO task inspired by the real-world application which are commonly used, based on gym library and can be widely used for evaluating MARL algorithm before it is implemented to the real application.

## 5.1 Earth Observation Satellite Mission

In Stephenson and Schaub (2024b), an open-source Python package for developing and customizing reinforcement learning environment has been designed by the authors to solve spacecraft tasking problems. It integrates Basilisk, a high-performance and high-fidelity spacecraft simulation framework, with the abstractions of satellite tasks and operational goals, all within the standard Gymnasium API wrapper for RL environments (see Fig. 5). This package is specifically designed to support the needs of researchers in reinforcement learning and spacecraft operations. Some works has been done to solve earth observation mission using RL algorithm Herrmann et al. (2023),Herrmann et al. (2024),Stephenson and Schaub (2024a) and currently a MARL framework has been provided as an example. This work has an open source codes available here: https://github.com/AVSLab/bsk_rl.

## 5.2 Power Grid Networks

A common power management benchmark has been proposed in Wang et al. (2021). This environment introduces a power network problem that provides a compelling yet challenging real-world scenario for the application of MARL as illustrated in Fig 6. The growing trend of de-carbonization is putting significant pressure on power distribution networks. Active voltage control has emerged as a promising solution to alleviate power congestion and enhance voltage stability without requiring additional hardware, by leveraging controllable devices within the network, such as rooftop Photo-voltaic (PV) and Static Var Compensator (SVC). It has been used for a robust evaluation of MARL algorithms in Guo et al. (2022), cooperative MARL with individual global max Hong et al. (2022), dynamic MARL algorithm configuration Xue et al. (2022), some power networks research Chen et al. (2022); Lu et al. (2023), etc. These devices are numerous and spread across wide geographic areas, making MARL an ideal approach. The codes of this work is publicly available here: https://github.com/Future-Power-Networks/MAPDN.

## 5.3 Traffic Management

CityFlow is a multi-agent reinforcement learning (MARL) environment developed for large-scale urban traffic management (Zhang et al., 2019). It tackles the complex task of traffic signal control, which requires real-time adaptation to dynamic traffic conditions and coordination among thousands of agents, such as vehicles and pedestrians. CityFlow features fundamentally optimized data structures and efficient algorithms,
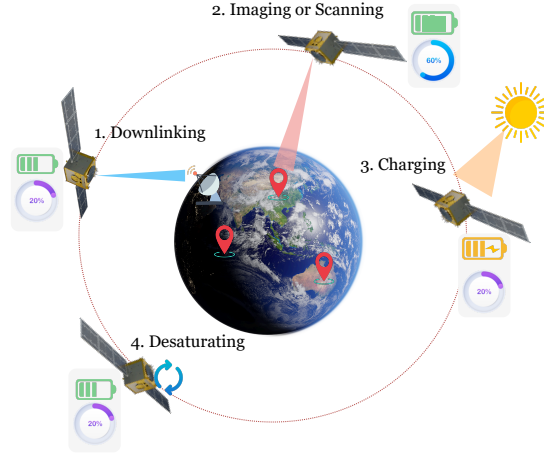
**Fig. 5**: Earth Observation Mission Environment in BSK-RL: The number of satellite can be defined either single satellite or multi-satellite. Each satellite have each actions: 1) Downlinking: Transmit collected data to the ground station with predefined transmission speed and delete the data after successfully downlinked; 2) Imaging or Scanning: Satellite capturing image of a target using optical sensor or scanning any object on the earth surface using radar sensor (with two different payloads and can be used for different tasks); 3) Charging: The satellite is in charging mode and pointing its solar panel towards sun direction to maximize solar energy absorption; 4) Desaturating: There is a condition of the satellite's rotating wheel rotates saturatedly and the speed should be reduced to control their attitude.

enabling high-speed, city-wide simulations. It supports flexible road network and traffic flow configurations based on both synthetic and real-world data, and includes a user-friendly interface for integrating reinforcement learning models. Most importantly, CityFlow enables large-scale, interactive traffic simulations and opens new possibilities for advancing intelligent transportation systems through machine learning. In Wei et al. (2019) the CityFlow is used to develop the algorithm. The resources of this simulator can be found here: https://github.com/cityflow-project/CityFlow

## 5.4 Container Management

A real-world industrial control task inspired resource allocation environment has been proposed in Pendyala et al. (2024). The authors outline the real-world industrial control task that served as the basis for RL benchmark. This task stems from the final phase of a waste sorting process (see Fig. 8). The environment comprises a solid material transformation facility containing multiple containers and a much smaller number of Processing Units (PUs). These containers are continuously filled with material,
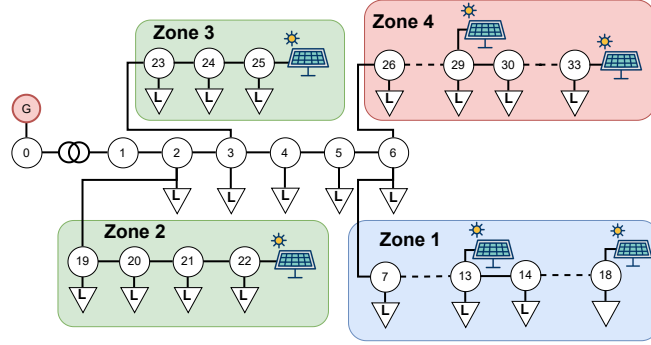
**Fig. 6**: Power Grid Networks environment: This system simulates the energy distribution and control within 4 different zones and 33 Bus Networks. Bus 2-33 voltages should be controlled by the system and Bus 0-1 represent the substation at the main grid with constant voltage and infinite active and reactive power capacity. There are several PV energy sources located in different areas that is challenging to maintain the voltage stability and control.
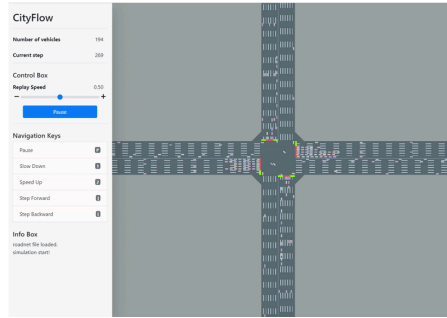


**Fig. 7**: Snapshot from the CityFlow simulation which manages traffic flows allocation at the intersection for single agent RL and multiple intersection for MARL

where the flow rate of the material follows a stochastic process that varies by container. Although this environment is developed for RL algorithm, it can be extended into the MARL settings by adding more agents. This work has a publicly available codes here: https://github.com/Pendu/ContainerGym.

### Tested Algorithms

Based on the available simulators, several algorithms have been found in the repository and reported in the literature as seen in Table 4.
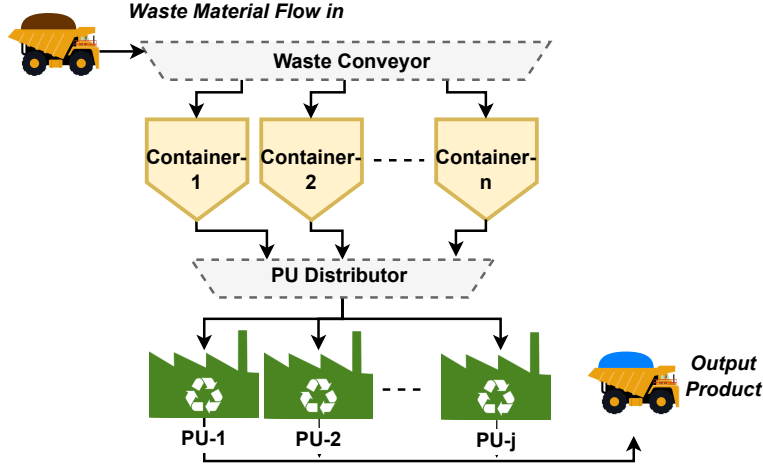
**Fig. 8**: ContainerGym environment: This environment simulate the waste processing unit that is assumed the waste materials flow input continuously and store them to multiple containers. It should maintain the container level and allocate to Processing Unit (PU) that requiring waste materials to be processed.

# 6 Future Directions and Potential Challenges

## 6.1 Future Directions

The application of MARL in RAO holds immense promise, with the potential to revolutionize numerous industries. As industrial systems become increasingly complex, decentralized, and dynamic, traditional resource allocation methods struggle to keep up. MARL's capabilities, such as distributed decision-making, real-time adaptation, and modeling intricate agent interactions, position it as a key enabler for future resource management systems.

A primary area of future research is enhancing scalability and efficiency. With the growing scale of systems, such as the Internet of Things (IoT), edge computing, and smart cities, MARL must handle more agents and larger environments. Techniques like hierarchical MARL, mean-field approximations Yang et al. (2018), and decentralized learning frameworks Zhang et al. (2019) will be critical in addressing computational complexity and communication overhead in such large-scale systems.

Dynamic and autonomous systems will also benefit from MARL advancements. For instance, next-generation wireless networks (e.g., 6G) Du et al. (2022), energy-efficient resource management Shen et al. (2022), and autonomous transportation systems can leverage MARL to enable real-time and adaptive resource allocation. These applications demand decentralized learning capabilities and the ability to optimize system performance while ensuring sustainability.

Inter-agent coordination and communication remain vital research areas. As systems grow to be more complex, seamless collaboration between agents becomes

35

**Table 4**: Summary of algorithms tested in the benchmarks or simulators

| Benchmark | Algorithm |
|---|---|
| **BSK-RL** | PPO |
| | MAPPO + Communication |
| **Power Distribution Networks** | IAC |
| | IDDPG |
| | MADDPG |
| | SQDDPG |
| | IPPO |
| | MAPPO |
| | MAAC |
| | MATD3 |
| | COMA |
| | FACMAC |
| **ContainerGym** | PPO |
| | TRPO |
| | DQN |
| **Traffic Management** | Graph based MARL |

increasingly important. Advances in coordination mechanisms, such as graph-based models Zhang et al. (2023) and communication-free methods, will improve resource sharing in cooperative environments while ensuring fairness in competitive settings.

Additionally, real-world resource allocation often involves non-stationary environments and heterogeneous agents Zhong et al. (2024). Future MARL research will focus on enhancing agent adaptability to evolving conditions and diverse objectives. By improving generalization across environments, MARL systems can rapidly adapt to new scenarios, making them more practical for dynamic resource allocation challenges.

Finally, integrating MARL with emerging technologies like federated learning Zhang et al. (2021); Li et al. (2022) and quantum computing Yun et al. (2023) presents exciting opportunities. These integrations could enable secure, efficient, and decentralized resource management systems, transforming domains like energy, logistics, and beyond by ensuring both efficiency and resilience in the face of growing complexity.

## 6.2 Potential Future Challenges

MARL offers a powerful solution for addressing the complexity and dynamic nature of RAO. By allowing decentralized agents to learn optimal strategies through interaction with their environment, MARL can handle the challenges of non-stationary condition, competition, and cooperation, providing a flexible, scalable approach for optimizing resource distribution in a variety of domains. However, several challenges remains in this domain such as: (1) Partial Observability Baker et al. (2020): Agents often have limited information, requiring approaches that can handle partial observability. (2) Convergence Yu et al. (2022): Careful algorithm design and hyperparameter tuning are needed to ensure convergence in competitive scenarios. (3) Scalability Yang et al. (2018); Wang et al. (2020): Techniques like mean field approximations are used to handle large numbers of agents. (4) Safety Constraints Lu et al. (2021): Implementing

safe exploration techniques to avoid harmful resource allocations during learning. (5) Communication Overhead Zhu et al. (2024): Centralized approaches require a high level of communication between agents and the central controller. This can lead to bottlenecks, especially in real-time environments or large-scale systems. (6) Exploration Cui et al. (2019): The joint exploration of the state-action space by all agents is more complex than in decentralized methods.

# 7 Conclusion

This survey examines the intersection of MARL and RAO, highlighting recent advances and emerging trends in this rapidly evolving area. MARL has demonstrated considerable potential in addressing the challenges of decentralized and dynamic resource allocation by enabling agents to learn and adapt in complex, uncertain environments. We reviewed core methodologies, surveyed applications across diverse domains, and analyzed the strengths and limitations of existing approaches.

Despite this progress, several key challenges remain, including non-stationary, limited scalability, coordination complexity, and the need for more generalizable algorithms. Tackling these issues calls for further research into improved training paradigms, adaptive communication mechanisms, and hybrid techniques that integrate MARL with classical optimization strategies.

Future work should also prioritize the deployment of MARL in emerging domains, the establishment of standardized RAO benchmarks, and the development of task-relevant evaluation metrics. By addressing these gaps, the community can advance the capabilities of MARL and broaden its impact on real-world resource allocation systems.

# References

Abid, M.S., H.J. Apon, S. Hossain, A. Ahmed, R. Ahshan, and M.H. Lipu. 2024. A novel multi-objective optimization based multi-agent deep reinforcement learning approach for microgrid resources planning. *Applied Energy* 353: 122029 .

Ahmed, M., J. Liu, M.A. Mirza, W.U. Khan, and F.N. Al-Wesabi. 2023. Marl based resource allocation scheme leveraging vehicular cloudlet in automotive-industry 5.0. *Journal of King Saud University-Computer and Information Sciences* 35 (6): 101420 .

Ahrarinouri, M., M. Rastegar, and A.R. Seifi. 2020. Multiagent reinforcement learning for energy management in residential buildings. *IEEE Transactions on Industrial Informatics* 17 (1): 659–666 .

Alam, M.R., M. St-Hilaire, and T. Kunz. 2016. Computational methods for residential energy cost optimization in smart grids: A survey. *ACM Computing Surveys (CSUR)* 49 (1): 1–34 .

Alcaraz, J. and C. Maroto. 2001. A robust genetic algorithm for resource allocation in project scheduling. *Annals of operations Research* 102: 83–109 .

Allahham, M.S., A.A. Abdellatif, N. Mhaisen, A. Mohamed, A. Erbad, and M. Guizani. 2022. Multi-agent reinforcement learning for network selection and resource allocation in heterogeneous multi-rat networks. *IEEE Transactions on Cognitive Communications and Networking 8*(2): 1287–1300 .

Antonio, G.P. and C. Maria-Dolores. 2022. Multi-agent deep reinforcement learning to manage connected autonomous vehicles at tomorrow's intersections. *IEEE Transactions on Vehicular Technology 71*(7): 7033–7043 .

Attiya, G. and Y. Hamam. 2006. Task allocation for maximizing reliability of distributed systems: A simulated annealing approach. *Journal of parallel and Distributed Computing 66*(10): 1259–1266 .

Baker, B., I. Kanitscheider, T. Markov, Y. Wu, G. Powell, B. McGrew, and I. Mordatch 2020. Emergent tool use from multi-agent autocurricula. In *International Conference on Learning Representations*.

Bi, J., H. Yuan, S. Duanmu, M. Zhou, and A. Abusorrah. 2020. Energy-optimized partial computation offloading in mobile-edge computing with genetic simulated-annealing-based particle swarm optimization. *IEEE Internet of Things Journal 8*(5): 3774–3785 .

Bratton, D. and J. Kennedy 2007. Defining a standard for particle swarm optimization. In *2007 IEEE swarm intelligence symposium*, pp. 120–127. IEEE.

Bu, L., R. Babu, B. De Schutter, et al. 2008. A comprehensive survey of multiagent reinforcement learning. *IEEE Trans. Syst., Man, Cybern. C, Appl., Rev. 38*(2): 156–172 .

Cao, Z., P. Zhou, R. Li, S. Huang, and D. Wu. 2020. Multiagent deep reinforcement learning for joint multichannel access and task offloading of mobile-edge computing in industry 4.0. *IEEE Internet of Things Journal 7*(7): 6201–6213 .

Cardon, A., T. Galinho, and J.P. Vacher. 2000. Genetic algorithms using multi-objectives in a multi-agent system. *Robotics and Autonomous systems 33*(2-3): 179–190 .

Cesana, M., I. Malanchini, and A. Capone 2008. Modelling network selection and resource allocation in wireless access networks with non-cooperative games. In *2008 5th IEEE International Conference on Mobile Ad Hoc and Sensor Systems*, pp. 404–409. IEEE.

Charbonnier, F., T. Morstyn, and M.D. McCulloch. 2022. Scalable multi-agent reinforcement learning for distributed control of residential energy flexibility. *Applied Energy* 314: 118825 .

Chen, D., M.R. Hajidavalloo, Z. Li, K. Chen, Y. Wang, L. Jiang, and Y. Wang. 2023. Deep multi-agent reinforcement learning for highway on-ramp merging in mixed traffic. *IEEE Transactions on Intelligent Transportation Systems 24*(11): 11623–11638 .

Chen, M., D. Gündüz, K. Huang, W. Saad, M. Bennis, A.V. Feljan, and H.V. Poor. 2021. Distributed learning in wireless networks: Recent progress and future challenges. *IEEE Journal on Selected Areas in Communications 39*(12): 3579–3605 .

Chen, P., S. Liu, X. Wang, and I. Kamwa. 2022. Physics-shielded multi-agent deep reinforcement learning for safe active voltage control with photovoltaic/battery energy

storage systems. *IEEE Transactions on Smart Grid 14*(4): 2656–2667 .

Chen, S., J. Dong, P. Ha, Y. Li, and S. Labi. 2021. Graph neural network and reinforcement learning for multi-agent cooperative control of connected autonomous vehicles. *Computer-Aided Civil and Infrastructure Engineering 36*(7): 838–857 .

Costa, B., L. Carvalho, M. Rosa, A. Araujo, et al. 2022. Computational resource allocation in fog computing: A comprehensive survey. *ACM Computing Surveys* .

Cui, J., Y. Liu, and A. Nallanathan. 2019. Multi-agent reinforcement learning-based resource allocation for uav networks. *IEEE Transactions on Wireless Communications 19*(2): 729–743 .

Du, X., T. Wang, Q. Feng, C. Ye, T. Tao, L. Wang, Y. Shi, and M. Chen. 2022. Multi-agent reinforcement learning for dynamic resource management in 6g in-x subnetworks. *IEEE transactions on wireless communications 22*(3): 1900–1914 .

Feriani, A. and E. Hossain. 2021. Single and multi-agent deep reinforcement learning for ai-enabled wireless networks: A tutorial. *IEEE Communications Surveys & Tutorials 23*(2): 1226–1252 .

Gao, H., X. Wang, W. Wei, A. Al-Dulaimi, and Y. Xu. 2023. Com-ddpg: Task offloading based on multiagent reinforcement learning for information-communication-enhanced mobile edge computing in the internet of vehicles. *IEEE Transactions on Vehicular Technology* .

Gao, X., R. Liu, and A. Kaushik. 2020. Hierarchical multi-agent optimization for resource allocation in cloud computing. *IEEE Transactions on Parallel and Distributed Systems 32*(3): 692–707 .

Gao, Z., L. Yang, and Y. Dai. 2022a. Fast adaptive task offloading and resource allocation via multiagent reinforcement learning in heterogeneous vehicular fog computing. *IEEE Internet of Things Journal 10*(8): 6818–6835 .

Gao, Z., L. Yang, and Y. Dai. 2022b. Large-scale computation offloading using a multi-agent reinforcement learning in heterogeneous multi-access edge computing. *IEEE Transactions on Mobile Computing 22*(6): 3425–3443 .

Gao, Z., L. Yang, and Y. Dai. 2023. Large-scale cooperative task offloading and resource allocation in heterogeneous mec systems via multi-agent reinforcement learning. *IEEE Internet of Things Journal* .

Gong, Y.J., J. Zhang, H.S.H. Chung, W.N. Chen, Z.H. Zhan, Y. Li, and Y.H. Shi. 2012. An efficient resource allocation scheme using particle swarm optimization. *IEEE Transactions on Evolutionary Computation 16*(6): 801–816 .

Guindo, L.A., M. Wagner, R. Baltussen, D. Rindress, J. van Til, P. Kind, and M.M. Goetghebeur. 2012. From efficacy to equity: Literature review of decision criteria for resource allocation and healthcare decisionmaking. *Cost Effectiveness and Resource Allocation 10*(1): 9. https://doi.org/10.1186/1478-7547-10-9 .

Guo, D., L. Tang, X. Zhang, and Y.C. Liang. 2020. Joint optimization of handover control and power allocation based on multi-agent deep reinforcement learning. *IEEE Transactions on Vehicular Technology 69*(11): 13124–13138 .

Guo, J., Y. Chen, Y. Hao, Z. Yin, Y. Yu, and S. Li 2022. Towards comprehensive testing on the robustness of cooperative multi-agent reinforcement learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 115–122.

Halabian, H. 2019a. Distributed resource allocation optimization in 5g virtualized networks. *IEEE Journal on Selected Areas in Communications 37*(3): 627–642 .

Halabian, H. 2019b. Distributed resource allocation optimization in 5g virtualized networks. *IEEE Journal on Selected Areas in Communications 37*(3): 627–642. https://doi.org/10.1109/JSAC.2019.2894305 .

Hao, J., T. Yang, H. Tang, C. Bai, J. Liu, Z. Meng, P. Liu, and Z. Wang. 2023. Exploration in deep reinforcement learning: From single-agent to multiagent domain. *IEEE Transactions on Neural Networks and Learning Systems* .

Heik, D., F. Bahrpeyma, and D. Reichelt. 2024. Adaptive manufacturing: dynamic resource allocation using multi-agent reinforcement learning .

Herrmann, A., M. Stephenson, and H. Schaub 2023. Reinforcement learning for multi-satellite agile earth observing scheduling under various communication assumptions. In *AAS Rocky Mountain GN&C Conference.*

Herrmann, A., M.A. Stephenson, and H. Schaub. 2024. Single-agent reinforcement learning for scalable earth-observing satellite constellation operations. *Journal of Spacecraft and Rockets 61*(1): 114–132 .

Hong, Y., Y. Jin, and Y. Tang. 2022. Rethinking individual global max in cooperative multi-agent reinforcement learning. *Advances in neural information processing systems* 35: 32438–32449 .

Hu, D., Z. Li, Z. Ye, Y. Peng, W. Xi, and T. Cai. 2024. Multi-agent graph reinforcement learning for decentralized volt-var control in power distribution systems. *International Journal of Electrical Power & Energy Systems* 155: 109531 .

Hu, J. and M.P. Wellman. 2003. Nash q-learning for general-sum stochastic games. *Journal of machine learning research 4*(Nov): 1039–1069 .

Hu, J., H. Zhang, L. Song, R. Schober, and H.V. Poor. 2020. Cooperative internet of uavs: Distributed trajectory design by multi-agent deep reinforcement learning. *IEEE Transactions on Communications 68*(11): 6807–6821 .

Huang, B., M. Zhou, X.S. Lu, and A. Abusorrah. 2023. Scheduling of resource allocation systems with timed petri nets: A survey. *ACM Computing Surveys 55*(11): 1–27 .

Ibaraki, T. and N. Katoh. 1988. *Resource allocation problems: algorithmic approaches.* MIT press.

Jain, V. and B. Kumar. 2023. Qos-aware task offloading in fog environment using multi-agent deep reinforcement learning. *Journal of Network and Systems Management 31*(1): 7 .

Jayanetti, A., S. Halgamuge, and R. Buyya. 2024. Multi-agent deep reinforcement learning framework for renewable energy-aware workflow scheduling on distributed cloud data centers. *IEEE Transactions on Parallel and Distributed Systems* .

Jendoubi, I. and F. Bouffard. 2023. Multi-agent hierarchical reinforcement learning for energy management. *Applied Energy* 332: 120500 .

Ji, Y., Y. Wang, H. Zhao, G. Gui, H. Gacanin, H. Sari, and F. Adachi. 2023. Multi-agent reinforcement learning resources allocation method using dueling double deep q-network in vehicular networks. *IEEE Transactions on Vehicular Technology 72*(10): 13447–13460 .

Jiandong, Z., Y. Qiming, S. Guoqing, L. Yi, and W. Yong. 2021. Uav cooperative air combat maneuver decision based on multi-agent reinforcement learning. *Journal of Systems Engineering and Electronics 32*(6): 1421–1438 .

Jiang, C. and Z. Sheng. 2009. Case-based reinforcement learning for dynamic inventory control in a multi-agent supply-chain system. *Expert Systems with Applications 36*(3): 6520–6526 .

Jiang, Y. 2015. A survey of task allocation and load balancing in distributed systems. *IEEE Transactions on Parallel and Distributed Systems 27*(2): 585–599 .

Jing, X., X. Yao, M. Liu, and J. Zhou. 2024. Multi-agent reinforcement learning based on graph convolutional network for flexible job shop scheduling. *Journal of Intelligent Manufacturing 35*(1): 75–93 .

Ju, Y., Y. Chen, Z. Cao, L. Liu, Q. Pei, M. Xiao, K. Ota, M. Dong, and V.C. Leung. 2023. Joint secure offloading and resource allocation for vehicular edge computing network: A multi-agent deep reinforcement learning approach. *IEEE Transactions on Intelligent Transportation Systems 24*(5): 5555–5569 .

Kang, H., X. Chang, J. Mišić, V.B. Mišić, J. Fan, and Y. Liu. 2023. Cooperative uav resource allocation and task offloading in hierarchical aerial computing systems: A mappo-based approach. *IEEE Internet of Things Journal 10*(12): 10497–10509 .

Kennedy, J. and R. Eberhart 1995. Particle swarm optimization. In *Proceedings of ICNN'95-international conference on neural networks*, Volume 4, pp. 1942–1948. ieee.

Khan, A.A., M. Abolhasan, W. Ni, J. Lipman, and A. Jamalipour. 2019. A hybrid-fuzzy logic guided genetic algorithm (h-flga) approach for resource optimization in 5g vanets. *IEEE Transactions on Vehicular Technology 68*(7): 6964–6974 .

Khan, S.U. and I. Ahmad 2006. Non-cooperative, semi-cooperative, and cooperative games-based grid resource allocation. In *Proceedings 20th IEEE International Parallel & Distributed Processing Symposium*, pp. 10–pp. IEEE.

Kim, Y. and H. Lim. 2021. Multi-agent reinforcement learning-based resource management for end-to-end network slicing. *IEEE Access* 9: 56178–56190. https://doi.org/10.1109/ACCESS.2021.3072435 .

Kirkpatrick, S., C.D. Gelatt Jr, and M.P. Vecchi. 1983. Optimization by simulated annealing. *science 220*(4598): 671–680 .

Konda, V. and J. Tsitsiklis. 1999. Actor-critic algorithms. *Advances in neural information processing systems* 12 .

Kosanoglu, F., M. Atmis, and H.H. Turan. 2024. A deep reinforcement learning assisted simulated annealing algorithm for a maintenance planning problem. *Annals of Operations Research 339*(1): 79–110 .

Kumari, A., R. Kakkar, S. Tanwar, D. Garg, Z. Polkowski, F. Alqahtani, and A. Tolba. 2024. Multi-agent-based decentralized residential energy management using deep reinforcement learning. *Journal of Building Engineering* 87: 109031 .

Lei, L., Y. Tan, K. Zheng, S. Liu, K. Zhang, and X. Shen. 2020. Deep reinforcement learning for autonomous internet of things: Model, applications and challenges. *IEEE Communications Surveys & Tutorials 22*(3): 1722–1760 .

Li, X., L. Lu, W. Ni, A. Jamalipour, D. Zhang, and H. Du. 2022. Federated multi-agent deep reinforcement learning for resource allocation of vehicle-to-vehicle

communications. *IEEE Transactions on Vehicular Technology 71* (8): 8810–8824 .

Liao, X., X. Hu, Z. Liu, S. Ma, L. Xu, X. Li, W. Wang, and F.M. Ghannouchi. 2020. Distributed intelligence: A verification for multi-agent drl-based multibeam satellite resource allocation. *IEEE Communications Letters 24* (12): 2785–2789 .

Lin, J.T. and C.C. Chiu. 2018. A hybrid particle swarm optimization with local search for stochastic resource allocation problem. *Journal of Intelligent Manufacturing 29* (3): 481–495 .

Liu, P., K. An, J. Lei, Y. Sun, W. Liu, and S. Chatzinotas. 2024. Computation rate maximization for scma-aided edge computing in iot networks: A multi-agent reinforcement learning approach. *IEEE Transactions on Wireless Communications* .

Liu, W., B. Li, W. Xie, Y. Dai, and Z. Fei. 2023. Energy efficient computation offloading in aerial edge networks with multi-agent cooperation. *IEEE Transactions on Wireless Communications 22* (9): 5725–5739 .

Liu, X.F., J. Zhang, and J. Wang. 2022. Cooperative particle swarm optimization with a bilevel resource allocation mechanism for large-scale dynamic optimization. *IEEE Transactions on Cybernetics 53* (2): 1000–1011 .

Liu, Y., J. Fan, L. Zhao, W. Shen, and C. Zhang. 2023. Integration of deep reinforcement learning and multi-agent system for dynamic scheduling of re-entrant hybrid flow shop considering worker fatigue and skill levels. *Robotics and Computer-Integrated Manufacturing* 84: 102605 .

Lowe, R., Y. Wu, A. Tamar, J. Harb, O.P. Abbeel, and I. Mordatch 2017. Multi-agent actor-critic for mixed cooperative-competitive environments. In *Advances in Neural Information Processing Systems*, pp. 6379–6390.

Lu, S., K. Zhang, T. Chen, T. Başar, and L. Horesh 2021. Decentralized policy gradient descent ascent for safe multi-agent reinforcement learning. In *Proceedings of the AAAI conference on artificial intelligence*, Volume 35, pp. 8767–8775.

Lu, Y., Y. Xiang, Y. Huang, B. Yu, L. Weng, and J. Liu. 2023. Deep reinforcement learning based optimal scheduling of active distribution system considering distributed generation, energy storage and flexible load. *Energy* 271: 127087 .

Ma, C., A. Li, Y. Du, H. Dong, and Y. Yang. 2024. Efficient and scalable reinforcement learning for large-scale network control. *Nature Machine Intelligence*: 1–15 .

Mao, H., M. Alizadeh, I. Menache, and S. Kandula 2016. Resource management with deep reinforcement learning. In *Proceedings of the 15th ACM workshop on hot topics in networks*, pp. 50–56.

Meng, F., P. Chen, L. Wu, and J. Cheng. 2020. Power allocation in multi-user cellular networks: Deep reinforcement learning approaches. *IEEE Transactions on Wireless Communications 19* (10): 6255–6267 .

Mnih, V., A.P. Badia, M. Mirza, A. Graves, T. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu 2016. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pp. 1928–1937. PMLR.

Naderializadeh, N., J.J. Sydir, M. Simsek, and H. Nikopour. 2021. Resource management in wireless networks via multi-agent deep reinforcement learning. *IEEE Transactions on Wireless Communications 20* (6): 3507–3523 .

Nasir, Y.S. and D. Guo. 2019. Multi-agent deep reinforcement learning for dynamic power allocation in wireless networks. *IEEE Journal on selected areas in communications 37*(10): 2239–2250 .

Nguyen, T.T., N.D. Nguyen, and S. Nahavandi. 2020. Deep reinforcement learning for multiagent systems: A review of challenges, solutions, and applications. *IEEE transactions on cybernetics 50*(9): 3826–3839 .

Ning, Z. and L. Xie. 2024. A survey on multi-agent reinforcement learning and its application. *Journal of Automation and Intelligence* .

Noor-A-Rahim, M., Z. Liu, H. Lee, G.M.N. Ali, D. Pesch, and P. Xiao. 2020. A survey on resource allocation in vehicular networks. *IEEE transactions on intelligent transportation systems 23*(2): 701–721 .

Oliehoek, F.A., C. Amato, et al. 2016. *A concise introduction to decentralized POMDPs*, Volume 1. Springer.

Orr, J. and A. Dutta. 2023. Multi-agent deep reinforcement learning for multi-robot applications: A survey. *Sensors 23*(7): 3625 .

Parvini, M., M.R. Javan, N. Mokari, B. Abbasi, and E.A. Jorswieck. 2023. Aoi-aware resource allocation for platoon-based c-v2x networks via multi-agent multi-task reinforcement learning. *IEEE Transactions on Vehicular Technology 72*(8): 9880–9896 .

Patriksson, M. 2008. A survey on the continuous nonlinear resource allocation problem. *European Journal of Operational Research 185*(1): 1–46 .

Pendyala, A., J. Dettmer, T. Glasmachers, and A. Atamna 2024. Containergym: A real-world reinforcement learning benchmark for resource allocation. In G. Nicosia, V. Ojha, E. La Malfa, G. La Malfa, P. M. Pardalos, and R. Umeton (Eds.), *Machine Learning, Optimization, and Data Science*, Cham, pp. 78–92. Springer Nature Switzerland.

Rashid, T., M. Samvelyan, C.S. De Witt, G. Farquhar, J. Foerster, and S. Whiteson. 2020. Monotonic value function factorisation for deep multi-agent reinforcement learning. *Journal of Machine Learning Research 21*(178): 1–51 .

Ren, L., X. Fan, J. Cui, Z. Shen, Y. Lv, and G. Xiong. 2022. A multi-agent reinforcement learning method with route recorders for vehicle routing in supply chain management. *IEEE Transactions on Intelligent Transportation Systems 23*(9): 16410–16420 .

Roesch, M., C. Linder, R. Zimmermann, A. Rudolf, A. Hohmann, and G. Reinhart. 2020. Smart grid for industry using multi-agent reinforcement learning. *Applied Sciences 10*(19): 6900 .

Saaty, T.L., L.G. Vargas, and K. Dellmann. 2003. The allocation of intangible resources: the analytic hierarchy process and linear programming. *Socio-Economic Planning Sciences 37*(3): 169–184 .

Sadatdiynov, K., L. Cui, L. Zhang, J.Z. Huang, S. Salloum, and M.S. Mahmud. 2023. A review of optimization methods for computation offloading in edge computing networks. *Digital Communications and Networks 9*(2): 450–461 .

Sarah, A., G. Nencioni, and M.M.I. Khan. 2023. Resource allocation in multi-access edge computing for 5g-and-beyond networks. *Computer Networks 227*: 109720 .

Schulman, J. 2015. Trust region policy optimization. *arXiv preprint arXiv:1502.05477* .

Schulman, J., F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347* .

Seid, A.M., G.O. Boateng, B. Mareri, G. Sun, and W. Jiang. 2021. Multi-agent drl for task offloading and resource allocation in multi-uav enabled iot edge network. *IEEE Transactions on Network and Service Management 18*(4): 4531–4547 .

Shao, X., F.S. Kshitij, and C.S. Kim. 2024. Gails: an effective multi-object job shop scheduler based on genetic algorithm and iterative local search. *Scientific Reports 14*(1): 2068 .

Shen, R., S. Zhong, X. Wen, Q. An, R. Zheng, Y. Li, and J. Zhao. 2022. Multi-agent deep reinforcement learning optimization framework for building energy system with renewable energy. *Applied Energy* 312: 118724 .

Singh, A.K., P. Dziurzanski, H.R. Mendis, and L.S. Indrusiak. 2017, April. A survey and comparative study of hard and soft real-time dynamic resource allocation strategies for multi-/many-core systems. *ACM Comput. Surv. 50*(2). https://doi.org/10.1145/3057267 .

Spinellis, D., C. Papadopoulos, and J.M. Smith. 2000. Large production line optimization using simulated annealing. *International journal of production research 38*(3): 509–541 .

Stephenson, M. and H. Schaub 2024a. Reinforcement learning for earth-observing satellite autonomy with event-based task intervals. In *AAS Rocky Mountain GN&C Conference, Breckenridge, CO.*

Stephenson, M.A. and H. Schaub 2024b. Bsk-rl: Modular, high-fidelity reinforcement learning environments for spacecraft tasking. In *75th International Astronautical Congress, Milan, Italy, IAF.*

Suman, B. and P. Kumar. 2006. A survey of simulated annealing as a tool for single and multiobjective optimization. *Journal of the operational research society 57*(10): 1143–1160 .

Sun, X. and J. Qiu. 2021. Two-stage volt/var control in active distribution networks with multi-agent deep reinforcement learning method. *IEEE Transactions on Smart Grid 12*(4): 2903–2912 .

Sunehag, P., G. Lever, A. Gruslys, W.M. Czarnecki, V. Zambaldi, M. Jaderberg, M. Lanctot, N. Sonnerat, J.Z. Leibo, K. Tuyls, and T. Graepel 2018. Value-decomposition networks for cooperative multi-agent learning based on team reward. In *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*, AAMAS '18, Richland, SC, pp. 2085–2087. International Foundation for Autonomous Agents and Multiagent Systems.

Sutton, R.S. 1988. Learning to predict by the methods of temporal differences. *Machine learning* 3: 9–44 .

Sutton, R.S. and A.G. Barto. 2018. *Reinforcement learning: An introduction.* MIT press.

Suzuki, A., R. Kawahara, and S. Harada. 2022. Cooperative multi-agent deep reinforcement learning for dynamic virtual network allocation with traffic fluctuations. *IEEE Transactions on Network and Service Management 19*(3): 1982–2000 .

Tang, J., D.K. So, E. Alsusa, K.A. Hamdi, and A. Shojaeifard. 2015. Resource allocation for energy efficiency optimization in heterogeneous networks. *IEEE Journal*

on *Selected Areas in Communications 33*(10): 2104–2117 .

Tseng, F.H., X. Wang, L.D. Chou, H.C. Chao, and V.C. Leung. 2017. Dynamic resource prediction and allocation for cloud data center using the multiobjective genetic algorithm. *IEEE Systems Journal 12*(2): 1688–1699 .

Ushakov, I.A. 2013. *Optimal resource allocation: with practical statistical applications and theory.* John Wiley & Sons.

Van Hasselt, H., A. Guez, and D. Silver 2016. Deep reinforcement learning with double q-learning. In *Proceedings of the AAAI conference on artificial intelligence*, Volume 30.

Vengerov, D. 2007. A reinforcement learning approach to dynamic resource allocation. *Engineering Applications of Artificial Intelligence 20*(3): 383–390 .

Wang, J., W. Xu, Y. Gu, W. Song, and T.C. Green 2021. Multi-agent reinforcement learning for active voltage control on power distribution networks. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan (Eds.), *Advances in Neural Information Processing Systems*, Volume 34, pp. 3271–3284. Curran Associates, Inc.

Wang, Q., W. Li, and A. Mohajer. 2024. Load-aware continuous-time optimization for multi-agent systems: Toward dynamic resource allocation and real-time adaptability. *Computer Networks* 250: 110526 .

Wang, X., L. Ke, Z. Qiao, and X. Chai. 2020. Large-scale traffic signal control using a novel multiagent reinforcement learning. *IEEE transactions on cybernetics 51*(1): 174–187 .

Wang, Y., T. Xu, X. Niu, C. Tan, E. Chen, and H. Xiong. 2020. Stmarl: A spatio-temporal multi-agent reinforcement learning approach for cooperative traffic light control. *IEEE Transactions on Mobile Computing 21*(6): 2228–2242 .

Wang, Z., T. Schaul, M. Hessel, H. Hasselt, M. Lanctot, and N. Freitas 2016. Dueling network architectures for deep reinforcement learning. In *International conference on machine learning*, pp. 1995–2003. PMLR.

Wei, H., C. Chen, G. Zheng, K. Wu, V. Gayah, K. Xu, and Z. Li 2019. Presslight: Learning max pressure control to coordinate traffic signals in arterial network. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, KDD '19, New York, NY, USA, pp. 1290–1298. Association for Computing Machinery.

Wei, W., R. Yang, H. Gu, W. Zhao, C. Chen, and S. Wan. 2021. Multi-objective optimization for resource allocation in vehicular cloud computing networks. *IEEE Transactions on Intelligent Transportation Systems 23*(12): 25536–25545 .

Wei, Z., B. Li, R. Zhang, X. Cheng, and L. Yang. 2023. Many-to-many task offloading in vehicular fog computing: A multi-agent deep reinforcement learning approach. *IEEE Transactions on Mobile Computing 23*(3): 2107–2122 .

Wen, G., J. Fu, P. Dai, and J. Zhou. 2021. Dtde: A new cooperative multi-agent reinforcement learning framework. *The Innovation 2*(4) .

Wong, A., T. Bäck, A.V. Kononova, and A. Plaat. 2023. Deep multiagent reinforcement learning: Challenges and directions. *Artificial Intelligence Review 56*(6): 5023–5056 .

Wu, G., Z. Xu, H. Zhang, S. Shen, and S. Yu. 2023. Multi-agent drl for joint completion delay and energy consumption with queuing theory in mec-based iiot. *Journal of Parallel and Distributed Computing* 176: 80–94 .

Wu, H., G.K.H. Pang, K.L. Choy, and H.Y. Lam. 2018. Dynamic resource allocation for parking lot electric vehicle recharging using heuristic fuzzy particle swarm optimization algorithm. *Applied Soft Computing* 71: 538–552 .

Wu, T., P. Zhou, K. Liu, Y. Yuan, X. Wang, H. Huang, and D.O. Wu. 2020. Multi-agent deep reinforcement learning for urban traffic light control in vehicular networks. *IEEE Transactions on Vehicular Technology 69*(8): 8243–8256 .

Xiao, Y., Y. Song, and J. Liu. 2023. Multi-agent deep reinforcement learning based resource allocation for ultra-reliable low-latency internet of controllable things. *IEEE Transactions on Wireless Communications 22*(8): 5414–5430 .

Xu, J., M. Zhao, J. Fortes, R. Carpenter, and M. Yousif. 2008. Autonomic resource management in virtualized data centers using fuzzy logic-based approaches. *Cluster Computing* 11: 213–227 .

Xu, X., Y. Jia, Y. Xu, Z. Xu, S. Chai, and C.S. Lai. 2020. A multi-agent reinforcement learning-based data-driven method for home energy management. *IEEE Transactions on Smart Grid 11*(4): 3201–3211 .

Xu, X., K. Xu, Z. Zeng, J. Tang, Y. He, G. Shi, and T. Zhang. 2024. Collaborative optimization of multi-energy multi-microgrid system: A hierarchical trust-region multi-agent reinforcement learning approach. *Applied Energy* 375: 123923 .

Xue, K., J. Xu, L. Yuan, M. Li, C. Qian, Z. Zhang, and Y. Yu. 2022. Multi-agent dynamic algorithm configuration. *Advances in Neural Information Processing Systems* 35: 20147–20161 .

Yang, Y., R. Luo, M. Li, M. Zhou, W. Zhang, and J. Wang 2018. Mean field multi-agent reinforcement learning. In *International conference on machine learning*, pp. 5571–5580. PMLR.

Ye, D. and J. Chen. 2013. Non-cooperative games on multidimensional resource allocation. *Future Generation Computer Systems 29*(6): 1345–1352 .

Yin, S. and F.R. Yu. 2021. Resource allocation and trajectory design in uav-aided cellular networks based on multiagent reinforcement learning. *IEEE Internet of Things Journal 9*(4): 2933–2943 .

You, X., X. Li, Y. Xu, H. Feng, J. Zhao, and H. Yan. 2020. Toward packet routing with fully distributed multiagent deep reinforcement learning. *IEEE Transactions on Systems, Man, and Cybernetics: Systems 52*(2): 855–868 .

Yu, C., A. Velu, E. Vinitsky, J. Gao, Y. Wang, A. Bayen, and Y. Wu. 2022. The surprising effectiveness of ppo in cooperative multi-agent games. *Advances in Neural Information Processing Systems* 35: 24611–24624 .

Yu, L., S. Qin, M. Zhang, C. Shen, T. Jiang, and X. Guan. 2021. A review of deep reinforcement learning for smart building energy management. *IEEE Internet of Things Journal 8*(15): 12046–12063 .

Yun, W.J., J.P. Kim, S. Jung, J.H. Kim, and J. Kim. 2023. Quantum multiagent actor–critic neural networks for internet-connected multirobot coordination in smart factory management. *IEEE Internet of Things Journal 10*(11): 9942–9952 .

Yánez, W., R. Mahmud, R. Bahsoon, Y. Zhang, and R. Buyya. 2020. Data allocation mechanism for internet-of-things systems with blockchain. *IEEE Internet of Things Journal 7*(4): 3509–3522. https://doi.org/10.1109/JIOT.2020.2972776 .

Zabihi, Z., A.M. Eftekhari Moghadam, and M.H. Rezvani. 2023, August. Reinforcement learning methods for computation offloading: A systematic review. *ACM Comput. Surv. 56*(1). https://doi.org/10.1145/3603703 .

Zeynivand, A., A. Javadpour, S. Bolouki, A.K. Sangaiah, F. Ja'fari, P. Pinto, and W. Zhang. 2022. Traffic flow control using multi-agent reinforcement learning. *Journal of Network and Computer Applications* 207: 103497 .

Zhang, B., W. Hu, A.M. Ghias, X. Xu, and Z. Chen. 2023. Multi-agent deep reinforcement learning based distributed control architecture for interconnected multi-energy microgrid energy management and optimization. *Energy Conversion and Management* 277: 116647 .

Zhang, G., K. Yang, and H.H. Chen. 2012. Resource allocation for wireless cooperative networks: A unified cooperative bargaining game theoretic framework. *IEEE Wireless Communications 19*(2): 38–43 .

Zhang, H., S. Feng, C. Liu, Y. Ding, Y. Zhu, Z. Zhou, W. Zhang, Y. Yu, H. Jin, and Z. Li 2019. Cityflow: A multi-agent reinforcement learning environment for large scale city traffic scenario. In *The World Wide Web Conference*, WWW '19, New York, NY, USA, pp. 3620–3624. Association for Computing Machinery.

Zhang, H., C. Jiang, N.C. Beaulieu, X. Chu, X. Wang, and T.Q. Quek. 2015. Resource allocation for cognitive small cell networks: A cooperative bargaining game theoretic approach. *IEEE Transactions on Wireless Communications 14*(6): 3481–3493 .

Zhang, H., H. Zhao, R. Liu, A. Kaushik, X. Gao, and S. Xu. 2024. Collaborative task offloading optimization for satellite mobile edge computing using multi-agent deep reinforcement learning. *IEEE Transactions on Vehicular Technology* .

Zhang, J.D., Z. He, W.H. Chan, and C.Y. Chow. 2023. Deepmag: Deep reinforcement learning with multi-agent graphs for flexible job shop scheduling. *Knowledge-Based Systems* 259: 110083 .

Zhang, K., J. Cao, and Y. Zhang. 2021. Adaptive digital twin and multiagent deep reinforcement learning for vehicular edge computing and networks. *IEEE Transactions on Industrial Informatics 18*(2): 1405–1413 .

Zhang, M., Y. Dou, P.H.J. Chong, H.C. Chan, and B.C. Seet. 2021. Fuzzy logic-based resource allocation algorithm for v2x communications in 5g cellular networks. *IEEE Journal on Selected Areas in Communications 39*(8): 2501–2513 .

Zhang, S.Q., Q. Zhang, and J. Lin. 2019. Efficient communication in multi-agent reinforcement learning via variance based control. *Advances in neural information processing systems* 32 .

Zhang, W., D. Yang, W. Wu, H. Peng, N. Zhang, H. Zhang, and X. Shen. 2021. Optimizing federated learning in distributed industrial iot: A multi-agent approach. *IEEE Journal on Selected Areas in Communications 39*(12): 3688–3703 .

Zhang, X. and S. Debroy. 2023. Resource management in mobile edge computing: a comprehensive survey. *ACM Computing Surveys 55*(13s): 1–37 .

Zhang, Y., Z. Mou, F. Gao, J. Jiang, R. Ding, and Z. Han. 2020. Uav-enabled secure communications by multi-agent deep reinforcement learning. *IEEE Transactions*

on *Vehicular Technology 69* (10): 11599–11611 .

Zhang, Y., G. Zheng, Z. Liu, Q. Li, and H. Zeng. 2024. Marlens: understanding multi-agent reinforcement learning for traffic signal control via visual analytics. *IEEE transactions on visualization and computer graphics* .

Zhang, Y., H. Zhu, D. Tang, T. Zhou, and Y. Gui. 2022. Dynamic job shop scheduling based on deep reinforcement learning for multi-agent manufacturing systems. *Robotics and Computer-Integrated Manufacturing* 78: 102412 .

Zhao, J., F. Hu, J. Li, and Y. Nie. 2023. Multi-agent deep reinforcement learning based resource management in heterogeneous v2x networks. *Digital Communications and Networks* .

Zhao, N., Z. Ye, Y. Pei, Y.C. Liang, and D. Niyato. 2022. Multi-agent deep reinforcement learning for task offloading in uav-assisted mobile edge computing. *IEEE Transactions on Wireless Communications 21* (9): 6949–6960 .

Zheng, P., L. Xia, C. Li, X. Li, and B. Liu. 2021. Towards self-x cognitive manufacturing network: An industrial knowledge graph-based multi-agent reinforcement learning approach. *Journal of Manufacturing Systems* 61: 16–26 .

Zhong, Y., J.G. Kuba, X. Feng, S. Hu, J. Ji, and Y. Yang. 2024. Heterogeneous-agent reinforcement learning. *Journal of Machine Learning Research* 25: 1–67 .

Zhou, H., K. Jiang, S. He, G. Min, and J. Wu. 2023. Distributed deep multi-agent reinforcement learning for cooperative edge caching in internet-of-vehicles. *IEEE Transactions on Wireless Communications 22* (12): 9595–9609 .

Zhu, C., M. Dastani, and S. Wang 2024. A survey of multi-agent deep reinforcement learning with communication. In *Proceedings of the 23rd International Conference on Autonomous Agents and Multiagent Systems*, AAMAS '24, Richland, SC, pp. 2845–2847. International Foundation for Autonomous Agents and Multiagent Systems.

Zhu, X., Y. Luo, A. Liu, M.Z.A. Bhuiyan, and S. Zhang. 2020. Multiagent deep reinforcement learning for vehicular computation offloading in iot. *IEEE Internet of Things Journal 8* (12): 9763–9773 .