

Age of Information Optimization in Laser-charged UAV-assisted IoT Networks: A Multi-agent Deep Reinforcement Learning Method

Geng Sun, Likun Zhang, Jiahui Li, Jing Wu, Jiacheng Wang, Zemin Sun, Changyuan Zhao,
Victor C.M. Leung, *Life Fellow, IEEE*

Abstract—The integration of unmanned aerial vehicles (UAVs) with Internet of Things (IoT) networks offers promising solutions for efficient data collection. However, the limited energy capacity of UAVs remains a significant challenge. In this case, laser beam directors (LBDs) have emerged as an effective technology for wireless charging of UAVs during operation, thereby enabling sustained data collection without frequent returns to charging stations (CSs). In this work, we investigate the age of information (AoI) optimization in LBD-powered UAV-assisted IoT networks, where multiple UAVs collect data from distributed IoTs while being recharged by laser beams. We formulate a joint optimization problem that aims to minimize the peak AoI while determining optimal UAV trajectories and laser charging strategies. This problem is particularly challenging due to its non-convex nature, complex temporal dependencies, and the need to balance data collection efficiency with energy consumption constraints. To address these challenges, we propose a novel multi-agent proximal policy optimization with temporal memory and multi-agent coordination (MAPPO-TM) framework. Specifically, MAPPO-TM incorporates temporal memory mechanisms to capture the dynamic nature of UAV operations and facilitates effective coordination among multiple UAVs through decentralized learning while considering global system objectives. Simulation results demonstrate that the proposed MAPPO-TM algorithm outperforms conventional approaches in terms of peak AoI minimization and energy efficiency. Ideally, the proposed algorithm achieves up to 15.1% reduction in peak AoI compared to conventional multi-agent deep reinforcement learning (MADRL) methods.

Index Terms—Age of information, laser-powered UAV systems, multi-agent reinforcement learning, and IoT data collection

This study is supported in part by the National Natural Science Foundation of China (62272194, 62471200), in part by the Science and Technology Development Plan Project of Jilin Province (20250101027JJ), in part by the Postdoctoral Fellowship Program of China Postdoctoral Science Foundation (GZC20240592), in part by China Postdoctoral Science Foundation General Fund (2024M761123), and in part by the Scientific Research Project of Jilin Provincial Department of Education (JJKH20250117KJ).

(Corresponding author: Jing Wu.)

Geng Sun is with the College of Computer Science and Technology, Key Laboratory of Symbolic Computation and Knowledge Engineering of Ministry of Education, Jilin University, Changchun 130012, China, and also with the College of Computing and Data Science, Nanyang Technological University, Singapore 639798 (e-mail: sungeng@jlu.edu.cn).

Likun Zhang, Jiahui Li, Jing Wu, and Zemin Sun are with the College of Computer Science and Technology, Jilin University, Changchun 130012, China (e-mails: zhanglk23@mails.jlu.edu.cn, lijiahui@jlu.edu.cn, wujing@jlu.edu.cn, sunzemin@jlu.edu.cn).

Jiacheng Wang and Changyuan Zhao are with the College of Computing and Data Science, Nanyang Technological University, Singapore 639798 (e-mails: jiacheng.wang@ntu.edu.sg, zhao0441@e.ntu.edu.sg).

Victor C. M. Leung is with the College of Computer Science and Software Engineering, Shenzhen University, Shenzhen 518060, China, and also with the Department of Electrical and Computer Engineering, University of British Columbia, Vancouver, BC V6T 1Z4, Canada (e-mail: vleung@ieee.org).

I. INTRODUCTION

The rapid development of unmanned aerial vehicles (UAVs) and the Internet of Things (IoTs) has sparked significant interest in their integration to create efficient and scalable systems for data collection, energy replenishment, and communication in various applications [1], [2]. This integration has emerged as a promising paradigm to address fundamental challenges in conventional IoT infrastructures, particularly regarding energy efficiency, communication range, and data collection capabilities. UAVs, characterized by their high mobility, flexibility, and ability to establish line-of-sight (LoS) communication links, offer a viable solution to the limitations of conventional IoT networks where battery life and communication range are often constrained [3].

However, one of the critical challenges in UAV-assisted IoT networks is the limited energy capacity of UAVs. In particular, the substantial propulsion energy consumption during flight operations, coupled with the energy requirements for communication and data processing, significantly restricts the operational duration of UAVs [4], [5]. Conventional approaches relying on fixed charging stations (CSs), while useful in some scenarios, introduce additional complexity and operational overhead, particularly when UAVs are deployed over large areas or in environments without accessible infrastructure [4], [6].

To overcome these limitations, recent advancements have proposed integrating UAVs with wireless charging technologies. Laser charging, in particular, has gained significant attention due to its ability to deliver concentrated energy over long distances [7], [8], thus allowing UAVs to recharge while still performing their data collection tasks [9]. Laser charging systems, such as those utilizing laser beam directors (LBDs) and high-power lasers, offer several advantages over traditional radio frequency (RF)-based charging, particularly in scenarios where UAVs must cover extensive areas or operate over long distances without returning to a fixed CS [10], [11]. However, deploying laser charging systems introduces its own set of challenges, including the need for accurate beam control [12], optimal energy transfer efficiency [13], and coordination between charging and data collection operations [14]. In particular, the performance metric of the data collection operations may have trade-offs with the charging efficiency of LBDs, which pose more requirements for designing the LBD-powered UAV systems.

In this work, we aim to investigate age of information (AoI) optimization in the LBD-powered UAV systems, which is a critical metric for evaluating the timeliness and freshness of information in IoT networks [5], [15]. This metric is particularly important in scenarios requiring real-time monitoring and decision-making, such as emergency response, industrial automation, and smart city applications. Minimizing AoI ensures that the collected data accurately represents the current state of the monitored environment, thereby enhancing the effectiveness of IoT-based systems [16], [17]. As aforementioned, the integration of AoI optimization with LBD-powered UAV systems presents a complex challenge that requires sophisticated approaches to balance information freshness with energy efficiency [8], [18].

Traditional optimization approaches such as convex optimization and mathematical programming have been widely applied to UAV trajectory planning and resource allocation problems [19]–[21]. While these methods provide optimal solutions under specific conditions, they often struggle with the real-time adaptability and dynamic nature required in LBD-powered UAV systems. Particularly, the precise beam tracking and continuous adjustment needed for effective laser charging demand solutions that can respond instantaneously to changing conditions, which conventional optimization techniques cannot adequately address [22]. In this case, deep reinforcement learning (DRL) approaches have gained significant attention due to their ability to adapt to dynamic environments and handle complex decision-making processes [23]. Various DRL techniques, including deep Q-network (DQN) [24], deep deterministic policy gradient (DDPG) [25], and proximal policy optimization (PPO) [26], have demonstrated promising results in optimizing UAV trajectories and resource allocation. However, these centralized DRL methods often fall short when dealing with the distributed characteristics of one-to-many charging scenarios common in LBD-powered UAV systems, where multiple UAVs need to coordinate their actions while operating in different regions [27].

Thus, in this work, we first conduct a comprehensive survey of the current state of LBD-powered UAV systems, identifying the key challenges and limitations of existing approaches. Based on this analysis, we seek to propose a novel distributed DRL framework that addresses the specific requirements of LBD-powered UAV-assisted IoT data collection systems, thereby enabling effective coordination among multiple UAVs while accounting for the temporal dependencies and spatial constraints inherent in these systems. The main contributions of this work are summarized as follows:

- *LBD-powered UAV Data Collection System:* We consider a UAV-assisted IoT network where UAVs collect data from IoTs while being recharged by laser beams from LBDs. In such systems, the UAVs operate in both charging and non-charging areas, collecting data from IoTs and optimizing the energy usage during the data collection process. This system is capable of providing efficient and sustainable data collection solutions in large-scale and energy-constrained environments [28].
- *Joint Optimization of UAV Trajectory and Laser Charging:* In the considered system, we formulate a joint

optimization problem that aims to minimize the peak AoI in the network by optimizing the UAV trajectories and charging strategies. This optimization actually needs to balance the data collection efficiency, energy consumption during flight and communication, and the effective use of laser charging for energy replenishment [29]. Thus, this problem is challenging due to its dynamic nature and the need to adapt to varying energy levels, IoT locations, and UAV positions in real-time.

- *MADRL-based Approach:* To solve the formulated optimization problem, we propose an MADRL-based approach, namely multi-agent proximal policy optimization with temporal memory and multi-agent coordination (MAPPO-TM). In MAPPO-TM, the incorporated temporal memory captures the dynamic nature of UAV trajectories and energy consumption while enabling coordination among multiple UAVs. This approach efficiently learns the optimal strategies through decentralized learning while considering the global system objectives.
- *Simulation and Performance Evaluation:* Simulation results demonstrate that the proposed MAPPO-TM algorithm outperforms conventional approaches in terms of peak AoI minimization and UAV energy efficiency. Thus, the algorithm effectively coordinates multiple UAVs, reduces the AoI, and optimizes the UAV trajectories and laser charging strategies.

The rest of this paper is organized as follows. Section II reviews the related works. Section III presents the models and problem formulation. Section IV introduces the proposed MAPPO-TM algorithm. In Section V, we present the simulation results. Finally, the paper is concluded in Section VII.

II. RELATED WORK

In this work, we aim to charge UAVs with lasers and optimize the flight trajectory of UAVs to better execute charging strategies and efficiently collect data in IoT networks. In the following, we will review some key related works to illustrate the novelty of our research. The comparisons between these related works and our work are shown in Table I.

A. Laser-powered UAVs Assisted Communications

1) *Laser-powered UAV System Architectures and Fundamentals:* Many existing works explored laser-powered UAV system architectures for data collection in IoT networks. For instance, the authors in [3] proposed a basic IoT structure integrating UAVs, IoTs, and LBDs, establishing the framework where UAVs collect data while LBDs provide energy. The authors in [10] explored the technical fundamentals of laser charging and classified optical wireless power transfer (OWPT) systems into laser power transfer (LPT) and LED-based OWPT, highlighting the advantages of LPT for high-density, long-distance energy transfer in applications like autonomous UAVs. Moreover, the authors in [11] categorized UAVs into charging UAVs (CUAVs) and mission UAVs (MUAVs) to enable aerial refueling without mission interruption. The authors in [30] developed a laser charging framework employing a drone-mounted base station (DBS)

TABLE I
COMPARISONS BETWEEN RELATED WORKS WITH THIS WORK.

Related works	Number of UAVs	Communication scenarios	Optimization variables	Network optimization metrics	Optimization methods
[30]	Single	UAV-aided relay communication for edge computing	DBS placement, bandwidth, power and UAV trajectory	DBS service time and task completion time for all UEs	Iterative algorithm and placement algorithm based on counting sort
[6]	Single	UAV-assisted MEC with HAP wireless charging	UAV flight trajectories	Energy efficiency and number of computation tasks collected by the UAV	Multi-objective reinforcement learning with trace-based experience replay
[9]	Single	Laser-powered UAV relay system for URLLC connectivity between source node and ground station	UAV trajectory, blocklength allocation, power control, and EH	Total decoding error rate, energy efficiency, and UAV energy consumption	Perturbation-based iterative method with divide-and-conquer approach
[29]	Single	Data gathering from ground IoTs with laser charging from HAPs	UAV hovering positions and charging energies at each position	Total task completion time (including data collection and charging time)	BCD approach, SCA and dynamic programming
[43]	Single	Air-ground coordinated MEC system with laser-powered UAV serving as both MEC server and relay for ground access point	Trajectory of UAV, computation task allocation between UAV and access point, and EH time allocation	Long-term average energy consumption of UAV	Two-step alternating optimization algorithm that combines linear programming for task and time allocation with DDPG for trajectory design
[46]	Multiple	UAV-assisted wireless powered MEC for metaverse applications	Charging time, computation tasks scheduling, UAV trajectory design	Computation efficiency	Two-stage alternating optimization algorithm based on multi-task DRL
[47]	Single	UAV-based data collection and wireless charging for IoTs	3D trajectory of UAV, IoT scheduling	Residual energy while satisfying IoT requirements	SCA and BCD algorithm
[5]	Single	UAV-enabled wireless sensor networks for data collection	Flight trajectory of UAV	AoI minimization and average AoI minimization	Dynamic programming and genetic algorithm
[17]	Multiple	Multi-UAV-assisted wireless backscatter networks for sensing data collection	Access control of ground users, beamforming and trajectory planning of UAVs	Long-term time-averaged AoI	Lyapunov optimization, BCD, soft actor-critic algorithm
[56]	Single	Aerial-ground collaborative MEC	UAV flight paths, task offloading ratios	Total AoI of ground devices and total energy consumption of UAV	Multi-objective learning algorithm based on PPO
[57]	Single	Visible light communication-based vehicles-to-everything communication with cluster-based architecture	UAV flight paths, task offloading ratios	Energy efficiency and AoI	MADRL with TD3
[60]	Single	UAV collecting data from SNs distributed across multiple islands	Transmit power of SNs, clustering of islands, UAV flight trajectory	Long-term average AoI	Clustering-based dynamic adjustment of the shortest path algorithm
[67]	Multiple	UAV-aided information gathering from multiple sources over unreliable wireless channels	Transmission scheduling decisions	Throughput maximization with per-source AoI guarantee	Confidence bound based oracle stationary randomized sampling algorithm
[79]	Multiple	UAV-assisted MEC system	Task offloading decisions, computation resource allocation, UAV trajectory control	Task completion delay, UAV energy consumption, number of offloaded tasks	JTORATC approach using distributed splitting, threshold rounding, Karush-Kuhn-Tucker method, and SCA
[23]	Single	UAV collects status update packets from distributed sensors	UAV flight trajectory, transmission scheduling	Weighted sum of AoI	DRL method
[90]	Multiple	NOMA-enabled multi-UAV collaborative caching network	Caching decision, 3D trajectory planning, power allocation, subchannel reusing	System content retrieving delay	MAPPO and matching-DRL solution
Our work	Multiple	UAVs collect data from IoTs while recharged by LBDs	UAV trajectories and laser charging strategies	Peak AoI minimization	MAPPO-TM algorithm

that concurrently delivers services to ground-based user equipments (UEs) while harvesting energy transmitted from a laser CS mounted on the macro base station (MBS). Within their system architecture, both the MBS and DBS incorporate computational servers, enabling UEs to offload computational tasks to either station, while the DBS is precisely positioned at optimal locations to enhance uplink communications and computing services for ground UEs during each operational time slot. In addition, the authors in [31] investigated various deployment strategies for such systems, analyzing coverage ca-

pabilities under different laser charging conditions and optical turbulence through stochastic geometry. The authors in [32] proposed a dynamic OWPT system with overhead facilities housing laser transmitters and tracking cameras for continuous charging of moving vehicles. The authors in [33] examined the integration of laser-beamed wireless power transfer into high-altitude platform (HAP)-aided multiaccess edge computing systems serving HAP-connected aerial user equipments. By discretizing the three-dimensional (3D) coverage space of the HAP, they established a sophisticated multitier tile grid-

based spatial structure that furnished aerial location options in the form of tile grids for effective laser charging of aerial user equipments. The authors in [34] formulated an energy-constrained UAV-aided mobile edge-cloud continuum framework wherein offloaded tasks from ground IoTs can be cooperatively executed by UAVs functioning as edge servers and cloud servers connected to a ground base station (GBS) that serves as an access point. In their particular framework implementation, UAVs powered by laser beams transmitted from the GBS subsequently provide wireless charging capabilities to IoTs. However, most existing studies did not fully address the integration of continuous charging areas with non-charging areas, which significantly limited their practical application in real-world scenarios where UAVs should operate across diverse operational zones.

2) *Resource Allocation and Trajectory Optimization in Laser-powered System:* Many existing works have focused on optimizing resource allocation and flight trajectories in laser-powered systems. For instance, the authors in [4] proposed a cost-aware UAV deployment strategy to ensure high-quality communication and energy links between UAVs and ground users. The authors in [6] investigated energy-efficient trajectory optimization with wireless charging for UAV-assisted mobile edge computing (MEC). Moreover, the authors in [9] studied joint resource allocation, trajectory design, and energy harvesting (EH) to achieve ultra-reliable and low-latency communication (URLLC) in laser-powered UAV relay systems. The authors in [14] proposed a multimodal charging system with a two-layer model considering both inductive wireless power transfer and laser power scheduling. The authors in [29] introduced the minimum completion time trajectory and charging optimization algorithm, optimizing hovering positions and charging energies using block coordinate descent (BCD). In addition, the authors in [35] presented solutions to optimize energy efficiency of UAV relaying in IoT systems through trajectory planning and bandwidth allocation. The authors in [36] developed an improved clustering algorithm to determine optimal visiting order and entry points for IoT clusters, maximizing system energy efficiency. These works have significantly advanced trajectory and resource optimization in laser-powered UAV systems. The authors in [37] examined the simultaneous lightwave information and power transfer scheme for laser-powered decode-and-forward UAV relays functioning within an optical wireless backhaul. Their primary objective centered on determining the optimal allocation of received beam energy across the decoding circuit, transmitting circuit, and rotor block of the relay to maximize quality-of-service metrics including achievable data rate, outage probability, and error probability. The authors in [38] developed a laser charging-enabled DBS framework where a ground-based laser CS continuously delivers energy to an aerial quadrotor DBS providing communication services to multiple users. They formulated a comprehensive optimization problem addressing joint power and bandwidth assignment along with laser charging-enabled DBS placement to simultaneously maximize both flight duration and communication data rates. The authors in [39] explored a laser-charged UAV relaying network wherein a rotary-wing UAV serves as a data relay between

a GBS and UE while simultaneously receiving power from a dedicated CS via laser beam transmission. They addressed the critical optimization challenge of minimizing CS power consumption while ensuring the minimum data requirements of UE, consequently developing an algorithm that jointly optimizes UAV trajectory and charging power allocation. However, the aforementioned studies typically employed conventional optimization methods that struggled with real-time adaptability required for dynamic environments, limiting their effectiveness in scenarios with rapidly changing conditions or incomplete environmental information.

3) *Various Applications of Laser-powered Systems:* Many existing studies have investigated specific applications of laser-powered UAVs in various domains. For example, the authors in [40] employed green energy-powered base stations with laser chargers to extend UAV uptime for wireless charging and data backhauling in wireless rechargeable sensor networks. The authors in [41] investigated maximizing harvested data in laser-powered UAV-supported IoT deployments, enabling battery-free IoTs to establish communication links via bistatic backscattering. Moreover, the authors in [42] and [43] explored UAV-assisted edge computing, with the latter proposing an air-ground coordinated MEC system where a laser-powered UAV served as both an MEC server and relay. The authors in [44] investigated multi-UAV systems with full-duplex mobile users under URLLC constraints. In addition, the authors in [45] presented a UAV-assisted multiuser network using laser charging and simultaneous wireless information and power transfer. The authors in [46] proposed a two-stage optimization algorithm based on multi-task DRL for laser charging in wireless-powered metaverse scenarios. The authors in [47] investigated a solar-powered UAV system where the UAV simultaneously collects data from IoTs on the ground and charges them utilizing laser charging technology. Their objective focused on maximizing the residual energy of UAV while fulfilling IoT requirements through joint optimization of the 3D trajectory of UAV and IoT scheduling protocols. The authors in [48] examined the modeling of data collection from backscatter nodes utilizing UAVs that maintain sustainable operations through wireless energy transfer from laser based CSs. The authors in [49] introduced an energy-efficient laser-charged UAV-enabled rechargeable wireless sensor network environment wherein UAVs, energized by laser beams transmitted from ground-based stations, deliver services, gather data, and transfer energy to sensor nodes (SNs). Their research formulated a sophisticated joint optimization problem encompassing power allocation, dynamic charging strategy, and path planning with the dual objectives of minimizing task completion time and SN death time. However, most of the above works did not adequately address the challenge of balancing data collection efficiency with energy management across charging and non-charging areas, limiting their comprehensive application in complex IoT networks.

Different from these existing works, this work uniquely designs a laser-powered UAV system that effectively integrates operations across both charging and non-charging zones, employs advanced DRL methods for adaptive trajectory optimization in dynamic environments, and comprehensively

addresses the challenge of balancing data collection efficiency with energy management in complex IoT networks.

B. Optimizations of AoI in IoT Networks

1) *UAV Trajectory Planning for AoI Optimization:* Many existing works have studied UAV path planning to minimize AoI in IoT networks. For instance, the authors in [5] studied the age-optimal trajectory planning problem in UAV-enabled IoT networks, designing optimal trajectories to minimize both the age of the oldest sensed information and the average AoI of all IoTs. The authors in [7] explored data collection in UAV-assisted IoT networks powered by harvested energy, aiming to minimize the mission total time while ensuring each IoT receives required energy and transfers its sensed data. Moreover, the authors in [8] studied average AoI optimization by optimizing UAV trajectory in energy recharging networks, where the UAV collects data from IoTs and is replenished by ground chargers. The authors in [15] investigated average AoI minimization based on UAV trajectory and time allocation for EH and data collection, where the UAV serves as both mobile data collector and charger for IoTs. The authors in [16] presented a review of UAV-aided data collection focusing on DRL approaches to minimize AoI. In addition, the authors in [17] formulated a multi-stage stochastic optimization to minimize long-term time-averaged AoI by jointly optimizing access control, beamforming, and trajectory planning for multiple UAVs. The authors in [50] constructed a space-air-ground integrated network with satellites, HAPs, UAVs, and terrestrial IoTs, minimizing system AoI through UAV trajectory design and network configuration. The authors in [51] developed an optimization framework aimed at minimizing the total AoI of data collected by UAVs from ground IoT networks. Recognizing that the total AoI depends critically on both UAV flight time and data collection duration at hovering points, they conducted joint optimization of hovering point selection and the sequential visiting order to these locations. The authors in [52] introduced a sophisticated DRL-based proactive UAV trajectory planning algorithm capable of autonomously adjusting flight policies in response to channel variations while balancing the trade-off between energy transmission and data collection, ultimately achieving optimal system-level AoI under dynamic channel conditions. The authors in [53] addressed the complexity of multiple UAV operations by formulating a joint multi-UAV trajectory planning and data collection problem as a mixed integer nonlinear programming model, with the dual objectives of minimizing both AoI and energy consumption. However, most existing studies did not fundamentally solve the energy issues associated with extended UAV operation, which impacted their ability to achieve sustained information freshness in long-duration missions.

2) *Trade-offs between Energy Efficiency and AoI:* Many existing works have addressed the crucial balance between energy efficiency and information freshness in networked systems. For example, the authors in [54] incorporated AoI as a metric to ensure information freshness and designed an AoI-aware energy efficiency resource allocation scheme for satellite-based IoT networks. The authors in [55] tackled

the problem of selecting the optimal number of connections that is both AoI-optimal and energy-efficient by introducing an energy efficiency-peak AoI ratio to enable a trade-off between AoI and energy consumption. Moreover, the authors in [56] studied the AoI and energy trade-off in an aerial-ground collaborative MEC system, formulating a multi-objective optimization problem to simultaneously minimize total AoI and energy consumption by optimizing flight paths and task offloading ratios. The authors in [57] explored energy efficiency and AoI awareness in a cluster-based visible light communication vehicles-to-everything system, evaluating the impact of vehicle numbers on energy efficiency and AoI. The authors in [58] studied AoI optimization in information-gathering wireless networks where sources are equipped with batteries harvesting ambient energy, analyzing how energy arrival patterns and transmission policies influence average AoI. In addition, the authors in [59] investigated average AoI optimization in wireless-powered networks with directional charging, proposing an AoI-aware periodical charging scheduling algorithm. The authors in [60] investigated the long-term average AoI-minimal problem in a UAV-assisted wireless-powered communication network spanning multiple islands. The authors in [61] examined the challenge of rechargeable-UAV-aided timely data collection in IoT networks, wherein the UAV gathers status updates from multiple sensors while maintaining its energy level above a required threshold through recharging. To balance information freshness against energy consumption, they developed a Markov decision process framework designed to minimize the weighted sum of average total AoI and average recharging price. The authors in [62] investigated the AoI and energy trade-off within a system utilizing a UAV for data collection across multiple IoT nodes. To thoroughly analyze the interplay between AoI and energy consumption, they conducted joint optimization of collection time, UAV trajectory, and time slot duration. The authors in [63] explored task offloading challenges in a UAV-aided wireless powered edge computing system, emphasizing information freshness enhancement while maintaining UAV energy safety. To achieve minimum average AoI, they proposed a comprehensive joint optimization approach encompassing ground device wireless charging power, UAV flight trajectory, and offloading decisions. However, although most of the above works made significant contributions to understanding energy-AoI trade-offs, they did not fundamentally solve the energy issues, which limited their effectiveness in scenarios requiring extended operation periods.

3) *AoI in Specialized Networks and Applications:* Many existing studies have explored AoI optimization in diverse application domains with specific requirements. For instance, the authors in [18] introduced a theoretical framework for optimizing second-order behaviors of wireless networks, making it well-suited for modeling AoI and timely-throughput. The authors in [64] proposed a learning-based robust resource allocation considering overlapping interference and AoI-sensitive service requirements for ultra-dense Industrial IoT networks. Moreover, the authors in [65] constructed a blockchain-based remote intelligent healthcare system aiming to minimize both AoI and energy consumption of medical data

transmission. The authors in [66] investigated AoI impact on task-oriented multicasting in multi-cell non-orthogonal multiple access (NOMA) networks. The authors in [67] tackled the AoI-guaranteed transmission scheduling problem as an AoI-guaranteed multi-armed bandit problem. The authors in [68] proposed an AoI-aware waveform design scheme for cooperative joint radar-communications systems. The authors in [69] proposed a reconfigurable intelligent surface-assisted Internet of vehicles network, minimizing AoI of vehicle-to-infrastructure links while prioritizing vehicle-to-vehicle payload transmission. The authors in [70] studied high-speed railway mobile networks to optimize sensor scheduling and transmit power, thereby minimizing average AoI. The authors in [71] pioneered AoI use in autonomous driving systems, showing how optimizing AoI simultaneously minimizes response time and maximizes throughput. In addition, the authors in [72] presented an analytical framework establishing a dual-action guideline for minimizing average AoI in random access networks. The authors in [73] analyzed the impact of user velocity on peak AoI distribution for ground and aerial users. The authors in [74] investigated a decentralized UAV-aided MEC system tailored for smart agricultural applications, where processing nodes utilize network function virtualization technology. They formulated a sophisticated methodology for efficient network function virtualization orchestration that concurrently minimizes critical performance metrics. The authors in [75] analyzed a UAV-assisted IoT ecosystem featuring multiple UAVs operating in a comprehensive cycle, i.e., launching from a central data center, gathering data from distributed ground SNs, distributing information to various users, and subsequently returning to the data center. However, in the aforementioned studies, the optimal strategy for AoI optimization varied in different application scenarios, and generic approaches often failed to address the unique characteristics and constraints of specific deployments.

Different from these existing works, this work uniquely leverages laser charging technology to fundamentally address the energy limitations that hamper extended UAV operations, enabling sustained information freshness over longer mission durations. Based on this, we propose an optimization framework specifically tailored for laser-powered scenarios, effectively balancing AoI minimization with energy efficiency considerations in the context of LBD-powered UAV systems.

C. Optimization Methods for Various UAV-assisted IoT Networks

1) *Static Optimization Methods*: Many existing works have applied static optimization techniques to UAV-assisted IoT networks. For example, the authors in [19] introduced a novel trajectory planning framework for quadrotors landing on aerial vehicle carriers, where they combined a quadrotor trajectory planning method based on lossless convexification theory with a sequential convex programming approach, enabling autonomous landing on both stationary and moving aerial vehicle carriers in 3D space. The authors in [76] employed the A* algorithm, a well-established path planning approach, as a component of their two-step optimization methodology for mobile

nest path planning. Moreover, the authors in [77] tackled the multi-UAV cooperative path planning challenge by formulating it as a constrained optimization problem and introducing the evolutionary state estimation-based multi-strategy jellyfish search algorithm to identify high-quality trajectories. The authors in [78] developed a constrained decomposition-based multi-objective evolution algorithm to address the formulated constrained multi-objective optimization problem, which incorporated dual objective functions focused on energy-efficient offloading and safe path planning for UAVs. In addition, the authors in [79] proposed a joint task offloading, computation resource allocation, and UAV trajectory control (JTORATC) approach, where the task offloading sub-problem was resolved using distributed splitting and threshold rounding methods, the computation resource allocation sub-problem was addressed through Karush-Kuhn-Tucker optimization, and the UAV trajectory control sub-problem was solved via successive convex approximation (SCA) techniques. The authors in [80] developed two sophisticated algorithms, variable particle swarm optimization and twin variable neighborhood particle swarm optimization, to jointly optimize power, bandwidth, and UAV trajectories with the objective of minimizing AoI. The authors in [81] established a comprehensive start-to-end strategy incorporating association and planning mechanisms to minimize AoIs of two SNs through a methodical iterative three-step process. Initially, they determine the locations of data collection points (CPs) at which the UAVs hover to collect data and establish the SN-CP association using a density-based clustering algorithm. Subsequently, they cluster the CPs to form CP clusters and establish the CP-UAV association. Finally, leveraging the outcomes from the previous steps, they optimize the flight trajectories of the UAVs through an improved ant colony optimization algorithm while accounting for limited endurance capability constraints. However, these static optimization methods struggled with the dynamic nature of UAV-IoT systems, computational complexity in large-scale scenarios, and difficulty in adapting to changing environmental conditions without complete system information.

2) *DRL-based Optimization Methods*: Many existing works have applied DRL techniques to optimize energy efficiency and flight trajectories in UAV-assisted IoT networks. For instance, the authors in [11] utilized a dual UAV system (CUAVs and MUAVs) and employed DRL to minimize mission completion time by optimizing charging schedules and travel paths. Moreover, the authors in [23] proposed a DRL-based algorithm for UAV-assisted data collection to determine the optimal flight trajectory of the UAV and transmission scheduling of ground IoTs. The authors in [25] developed a UAV-aided video transmission system based on MEC and implemented a DDPG algorithm to achieve continuous action control through policy iteration. The authors in [26] proposed a PPO agent to enhance energy efficiency while addressing far-near fairness in NOMA-UAV networks by simultaneously controlling UAV trajectory, transmit power, node association, and power allocation. The authors in [76] devised a two-step optimization method employing modified multi-step dueling double DQN and A* algorithms to minimize inspection time while maximizing energy efficiency. In addition, the

authors in [82] examined a wireless-powered communication network where a resource-constrained secondary node harvests energy from ambient RF signals, implementing DRL to jointly optimize EH time and transmit power. The authors in [83] investigated a UAV-assisted IoT network in which the UAV sequentially accesses IoTs, proposing a DRL algorithm for multi-objective collaborative optimization that generates optimal strategies based on device priorities and assigned weights. The authors in [84] developed generative artificial intelligence agents for model formulation and subsequently implemented a mixture of experts (MoE) approach to design transmission strategies. Specifically, they harnessed large language models to construct an interactive modeling paradigm and employed retrieval-augmented generation to extract satellite expert knowledge that underpins mathematical modeling. Subsequently, through the integration of expertise from multiple specialized components, they introduced an MoE-PPO approach to address the formulated problem. The authors in [85] investigated reconfigurable intelligent surface (RIS)-assisted simultaneous wireless information and power transfer networks with rate splitting multiple access. To address the non-convex problem comprising both discrete and continuous variables, they proposed a DRL-based approach utilizing the PPO framework. Unlike traditional optimization approaches that optimize beamforming vectors and phase shifts separately and alternatively, their proposed PPO-based approach optimizes all variables simultaneously in unison. The authors in [86] formulated a double Q-learning-based trajectory design framework that enables energy-constrained surveillance UAVs to determine the optimal sequence of firefighting UAVs to visit (optimal flying trajectory), thereby maximizing the number of informed firefighting UAVs while accounting for limited execution time constraints. However, these DRL approaches were predominantly single-agent, centralized solutions that were inherently unsuitable for distributed environments, limiting their scalability and adaptability in multi-UAV scenarios.

Moreover, many existing works have investigated MADRL frameworks to optimize coordination and resource allocation in complex UAV networks. For example, the authors in [87] examined the use of UAV-enabled flying energy sources, powered by ground laser chargers, to support a set of MUAVs through RF wireless charging. They determined the positions of these flying energy sources using a MADRL approach called the multi-agent deep deterministic policy gradient (MADDPG) method, balancing fairness and energy consumption optimization. Moreover, the authors in [88] proposed a cloud-based task processing framework with MADRL that generates near real-time task offloading decisions based on partially knowable future information. The authors in [89] developed a double-level DRL approach within a divide-and-conquer framework, where upper-level DRL manages task allocation while lower-level DRL handles route planning. The authors in [90] introduced a multi-agent proximal policy optimization (MAPPO) algorithm for UAV collaborative caching and a matching-DRL solution for trajectory planning, power allocation, and channel reusing. The authors in [91] formulated a MADRL-based joint optimization scheme for trajectory, power control, user association, and subcarrier allocation

in multi-UAV networks. The authors in [92] implemented MADDPG to optimize task scheduling and UAV trajectory in UAV-based MEC. In addition, the authors in [93] devised a centralized MADRL algorithm to adjust aerial base station trajectories based on link quality estimations. The authors in [94] introduced a multi-agent cooperative DQN algorithm with delayed reward to minimize the AoI in scenarios where only one UAV can be charged at a time. The authors in [95] implemented a MADRL method to optimally control the trajectories of both unmanned ground vehicles (UGVs) and UAVs, thereby jointly reducing their energy consumption, decreasing the AoI of IoTs, and ensuring timely charging of UAVs while preventing their failures. The authors in [96] utilized a multi-agent deep Q-network (MADQN) algorithm to jointly optimize UAV trajectories, EH, task scheduling, and data offloading with the objective of minimizing the AoI and enhancing energy efficiency. Through the MADQN algorithm, UAVs can identify optimal data collection and EH decisions to minimize their energy consumption and efficiently gather data from multiple SNs, resulting in reduced AoI and improved energy efficiency. However, while these multi-agent approaches achieved remarkable success in UAV coordination, they were not specifically tailored for the unique challenges presented by LBD-powered systems, limiting their effectiveness in such specialized environments.

Different from these existing works, this work designs a novel LBD-powered UAV data collection system and proposes a corresponding MADRL framework specifically tailored for energy-efficient information gathering in IoT networks. The considered approach addresses the limitations of static optimization methods through adaptive learning-based techniques, overcomes the scalability constraints of single-agent DRL approaches through multi-agent coordination, and specifically optimizes for the unique characteristics of laser-powered environments.

III. SYSTEM MODEL

In this section, we first introduce the overall architecture of the LBD-powered multi-UAV data collection system in IoT networks. Then, we detail the data transmission model from IoTs to UAVs, the laser charging model from LBDs to UAVs, the AoI model, and the propulsion energy consumption model of UAVs. Finally, we formulate the problem. Note that the main notations used in this paper are summarized in Table II.

A. System Overview

As illustrated in Fig. 1, we consider an LBD-powered multi-UAV data collection system in IoT networks. Specifically, the system comprises two primary components, which are a set of IoTs denoted as $\mathcal{S} = \{1, 2, \dots, N_S\}$ and a set of UAVs denoted as $\mathcal{U} = \{1, 2, \dots, N_U\}$. Each IoT $s_i \in \mathcal{S}$ is placed at a known fixed location. We also consider that each IoT stores V_i units of data, which are static and do not increase over time, and has an initial energy level E_s . Moreover, the UAVs are tasked with visiting these IoTs to collect the stored data. We consider that each UAV $u_j \in \mathcal{U}$ has a fixed initial energy E_u and operates at a constant speed v at a fixed altitude H . In addition, LBDs

TABLE II
MAIN NOTATIONS.

Notation used in system model		Notation used in reinforcement learning	
Notation	Definition	Notation	Definition
A	Rotor disc area	\mathcal{A}	Action space
d_0	Fuselage drag ratio	\mathcal{A}_t	Action space for all UAVs at time slot t
E_i	Initial energy of each IoT	a_t^i	Action of UAV i at time slot t
E_u	Initial energy of each UAV	f_{actor}	Policy head that outputs action probabilities
H	Flight altitude of UAVs	h_t^i	Hidden state that captures temporal dependencies
N_D, N_S, N_U	The number of LBDs, IoTs and UAVs	\mathcal{O}	Observation space
P_i	Transmit power of the i -th IoT	\mathcal{O}_t	Observation space for all UAVs at time slot t
P_L	Laser transmitting power	o_t^i	Local observation of UAV i at time slot t
P_α, P_β	Constants representing the blade profile power and induced power during hover	\mathcal{R}	Immediate reward
R_c	Radius of charging area	r_t	Comprehensive reward function at time slot t
t, T, \mathcal{T}	The index, the number, and the set of time slots	\mathcal{S}	State space
t_d	Time duration of time slot	\mathcal{S}_t	System state at time slot t
v	Flight speed of UAVs	V_{local}^i	Individual performance value of UAV i
v_0	Mean rotor induced velocity in hover	V_{global}	Overall system state value
v_{tip}	Tip speed of the rotor blade	\mathcal{O}	Observation space
V_i	The amount of data to be transmitted each time generated by each IoT	\mathcal{O}_t	Observation space for all UAVs at time slot t
W	Channel bandwidth	θ	Parameter of LSTM-based actor networks for decentralized execution
x_i^s, y_i^s, z_i^s	Coordinate of the i -th IoT	ϕ	Parameter of centralized critic network for global state evaluation
$x_j(t), y_j(t), z_j(t)$	Coordinate of the j -th UAV in time slot t	θ_{old}	Parameter of old actor network for stable policy updates
x_k^d, y_k^d, z_k^d	Coordinate of the k -th LBD	γ	discount factor
β_0	Channel power at the reference distance 1 meter	π_θ	Policy network
δ	Laser attenuation coefficient	ω_l, ω_g	Learnable weights that balance local and global objectives
η_{le}	Laser-to-electricity conversion efficiency		
ρ	Air density		
σ^2	Gaussian noise power at the UAVs		
ω	Rotor solidity		

are denoted as $\mathcal{D} = \{1, 2, \dots, N_D\}$. The recharge process occurs when the UAV is within the charging area defined by radius R_c around an LBD, as shown in Fig. 1. Consequently, UAVs are deployed to collect data while operating under energy limitations.

We consider a discrete-time system that evolves in time slots $\mathcal{T} \triangleq \{1, 2, \dots, T\}$, where the length of a time slot is equal to t_d seconds. In this system, the UAVs collect data from any IoT when within its communication range. Since the energy of a UAV is depleted during flight and data communication, to prevent UAVs from exhausting their energy, they can be recharged at designated LBDs located within the operational area. Furthermore, data collection is considered to be instantaneous when the UAV is in the vicinity of an IoT.

We define a Cartesian coordinate system to describe the locations of IoTs, UAVs, and LBDs. Specifically, the IoTs are fixed at $(x_i^s, y_i^s, 0)$, representing their position at ground level. Moreover, UAVs fly at a constant altitude H , and their coordinates are represented as $(x_j(t), y_j(t), H)$ at any given time t , with $(x_j(t), y_j(t))$ denoting the horizontal position. Furthermore, the LBDs are located at known fixed coordinates (x_k^d, y_k^d, z_k^d) , and the charging area is defined by a radius R_c centered on each LBD, as illustrated in Fig. 1. Therefore, the operational area includes regions designated for IoT communication, UAV operation, and charging.

To establish a tractable theoretical framework for this complex system while ensuring meaningful optimization in the system, we make the following key assumptions:

- The UAVs operate at a fixed altitude H with constant speed v , enabling predictable flight patterns and simplify-

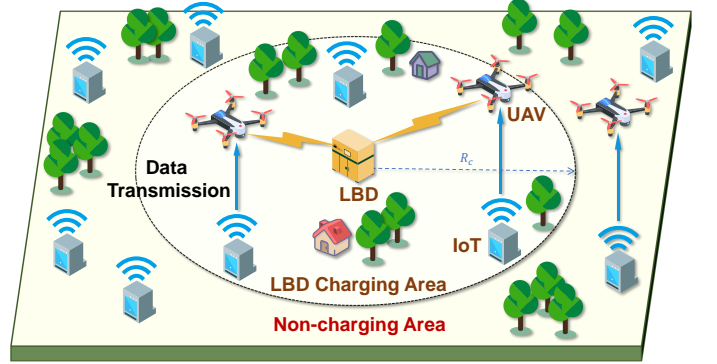


Fig. 1. Sketch map of the LBD-powered multi-UAV data collection system in IoT networks.

ing trajectory optimization. This assumption is commonly adopted in UAV trajectory optimization studies [97]–[100] as it reduces the complexity of the 3D path planning problem while remaining practical for many real-world UAV applications.

- All IoTs are stationary with fixed locations on the ground plane, each storing a static amount of data V_i that could be collected by UAVs. Note that this is reasonable for most IoT monitoring scenarios such as environmental monitoring, smart agriculture, and infrastructure inspection, where sensors are deployed at strategic locations to collect environmental or infrastructure data [101]–[104].

- *Data collection is instantaneous when a UAV is within communication range of an IoT, allowing for efficient modeling of the collection process.* This assumption is widely adopted in UAV-assisted IoT networks and is reasonable given the high mobility of UAVs and the relatively small data packet sizes typical in IoT applications [105]–[108].
- *LBDs are deployed at fixed locations, creating circular charging areas with radius R_c where UAVs can be recharged while continuing their operations.* This assumption follows established laser charging system models [109], [110] and reflects the practical deployment of ground-based laser charging infrastructure.
- *Laser charging from LBDs to UAVs is stable and reliable within the charging area, with charging power following the laser attenuation model without interruptions or fluctuations.* This assumption allows us to establish baseline performance and is supported by recent advancements in laser charging technology that have demonstrated high reliability in controlled environments [105], [111].

B. Data Transmission Model

Let P_i denote the transmit power of $s_i \in S$. Then, the data transmission rate from the IoT s_i to UAVs in bits/Hz at time slot $t \in \mathcal{T}$ is given by

$$R_{ij}^f(t) = W \log_2 \left(1 + \frac{\beta_0 P_i [P_{i,j}^{\text{LoS}}(t) \mu_{\text{LoS}} + P_{i,j}^{\text{NLoS}}(t) \mu_{\text{NLoS}}]}{\sigma^2 d_{ij}(t)^\alpha} \right), \quad (1)$$

where $d_{ij}(t) = \sqrt{(x_j(t) - x_i^s)^2 + (y_j(t) - y_i^s)^2}$ represents the distance between UAV j and IoT i at time t . Moreover, W represents the channel bandwidth, β_0 denotes the channel power at the reference distance of 1 meter, and σ^2 is the Gaussian noise power in UAVs. In addition, $P_{i,j}^{\text{LoS}}(t)$ and $P_{i,j}^{\text{NLoS}}(t)$ are the probability of LoS connection between sensor i and UAV j . Specifically, $P_{i,j}^{\text{LoS}}(t) = (1 + b_1 e^{-b_2(\theta_{i,j}(t) - b_1)})^{-1}$, and $P_{i,j}^{\text{NLoS}}(t) = 1 - P_{i,j}^{\text{LoS}}(t)$, in which b_1 and b_2 represent environment-dependent constants, and $\theta_{i,j}(t)$ is the angel between the IoT i and the UAV j .

C. Laser Charging Model

In the charging area, each UAV can be charged with one LBD. Let the radius of the charging area on a horizontal plane be R_c . Then, the power received from the UAV $u_j \in U$ from any LBD at time slot $t \in \mathcal{T}$ is given by

$$P_j^f(t) = P_L \cdot \eta_{le} \cdot e^{-\delta \cdot \sqrt{(x_j(t) - x^b)^2 + (y_j(t) - y^b)^2 + H^2}}, \quad (2)$$

where δ is the laser attenuation coefficient and η_{le} denotes the laser-to-electricity conversion efficiency. Moreover, P_L is the power of the laser.

D. AoI Model

The AoI depicts the time difference between the time an IoT generates data to be transmitted and the time data collection begins for this IoT. At first, all IoTs are equipped with

information and ready to be accessed by UAVs. The initial value of the instant of time is equal to 0. In time slot t , when a UAV accesses an IoT that has information to transmit, the AoI of IoT $a_i(t) = t - t'$, while t' represents the time when the IoT last generated information. Let $a_i(t)$ be the AoI of the IoT s_i at time instant t . When an IoT is accessed by a UAV and generates a larger AoI, the AoI of this IoT should be updated. Thus, the peak AoI of the network is given by

$$A = \max_{Q(t), T} \{a_1(t), a_2(t), \dots, a_n(t)\}. \quad (3)$$

E. Propulsion Energy Consumption Model of UAVs

According to dynamic principles, the propulsion power consumption of a UAV can be modeled as a function of its velocity. In this study, we adopt the propulsion power consumption model presented in [12], which is given by

$$P(v) = P_\alpha \left(1 + \frac{3v^2}{v_{\text{tip}}^2} \right) + P_\beta \left(\sqrt{1 + \frac{v^4}{4v_0^4} - \frac{v^2}{2v_0^2}} \right)^{\frac{1}{2}} + \frac{1}{2} d_0 \rho \omega A v^3, \quad (4)$$

where P_α and P_β are constants, with P_α representing the blade profile power and P_β representing the induced power during hover. Additionally, v_{tip} denotes the tip speed of the rotor blade, while v_0 is the mean rotor induced velocity in hover. The parameter d_0 corresponds to the fuselage drag ratio, and ω represents the rotor solidity. Moreover, ρ denotes the air density, and A is the area of the rotor disc.

From Eq. (4), it follows that the power consumption of the UAV at hover is $P(0) = P_\alpha + P_\beta$. The value of $P(0)$ is a constant that depends on the weight of the UAV, the air density, the rotor radius, and other factors. Based on Proposition 1 in [13], we conclude that Eq. (4) is convex for $v > 0$. Therefore, we can determine the optimal velocity $v_e = \arg \min P(v)$ that minimizes power consumption.

F. Problem Formulation

In this work, the primary optimization objective is to minimize the peak AoI across the LBD-powered multi-UAV data collection system in IoT networks. To achieve this goal, we need to carefully plan the flight trajectories of UAVs so that they can collect IoT information as efficiently and quickly as possible while adhering to charging strategies, thereby minimizing the peak AoI of the system.

To achieve the aforementioned optimization objective, we define the following key decision variables, including $\mathbf{A} = \{[a_u^x(t), a_u^y(t)] \mid t \in \mathcal{T}, u \in \mathcal{U}\}$ and $\mathbf{C} = \{c_u(t) \mid t \in \mathcal{T}, u \in \mathcal{U}\}$. Specifically, \mathbf{A} is a matrix representing the control parameters of the UAVs, which denotes its spatial displacement at different time slots. On the other hand, \mathbf{C} is a matrix representing the charging parameters of the UAVs, indicating whether the UAVs should leave or fly to the charging area, or continue to perform the missions at different time slots.

Following this, the optimization problem is formulated as follows:

$$\min_{A,C} \quad A = \max\{a_1(t), a_2(t), \dots, a_n(t)\}, \quad (5a)$$

$$\text{s.t.} \quad \int_0^T R_{ij}^f(t) dt \geq V_i, t \in \mathcal{T}, 1 \leq i \leq n, 1 \leq j \leq m, \quad (5b)$$

$$E_i - \int_0^T P_i(t) dt \geq E_\theta, t \in \mathcal{T}, 1 \leq i \leq n, \quad (5c)$$

$$0 < E_j(t) \leq E, t \in \mathcal{T}, 1 \leq j \leq m, \quad (5d)$$

$$\sqrt{(x_j(t) - x_{j'}(t))^2 + (y_j(t) - y_{j'}(t))^2} > d, \\ t \in \mathcal{T}, 1 \leq j \leq m, 1 \leq j' \leq m, j \neq j', \quad (5e)$$

$$\sqrt{(x_j(t) - x^b)^2 + (y_j(t) - y^b)^2} \leq 2R_c, \\ t \in \mathcal{T}, 1 \leq j \leq m, \quad (5f)$$

where constraint (5b) ensures that all data from IoTs can be collected by UAVs. Moreover, constraint (5c) guarantees that the remaining energy of s_i is greater than or equal to E_θ . In addition, constraint (5d) ensures that the remaining energy of each UAV is within a reasonable range. Furthermore, constraint (5e) is about collision control between the UAVs, and constraint (5f) limits the flight area.

As can be seen, the formulated optimization problem is NP-hard. To establish the computational complexity of our problem, we demonstrate its relationship to the well-known traveling salesman problem (TSP), which is NP-hard. Consider a simplified version of our problem with a single UAV, where energy constraints and collision avoidance are temporarily ignored. In this case, the problem reduces to finding the shortest path that visits each IoT exactly once, precisely the definition of TSP. Given that TSP is NP-hard, and our problem extends it by incorporating additional complexities such as multiple UAVs, energy constraints, charging strategies, and peak AoI minimization, we can conclude that our problem is also NP-hard.

IV. MADRL-BASED APPROACH

In this section, we propose an MADRL-based approach to solve our optimization problem. To this end, we first introduce the basics of DRL. Then, we show the motivations and rationales for using DRL and reformulate the problem as a partially observable Markov decision process (POMDP). Finally, we introduce the proposed MAPPO-TM algorithm with several improvements.

A. Overview of the basics of DRL

DRL represents a machine learning approach that integrates reinforcement learning with deep neural networks to address complex decision-making challenges in dynamic environments [112]. Central to DRL is the Markov decision process (MDP), which offers a mathematical framework for modeling sequential decision-making under uncertainty [113], [114]. Within an MDP, an agent interacts with an environment through discrete time steps, thereby making decisions that aim to maximize cumulative rewards. From a mathematical

perspective, an MDP comprises a tuple (S, A, P, R, γ) , where S denotes the state space, A represents the action space, P indicates the state transition probability function, R signifies the reward function, and γ is the discount factor. In this framework, the DRL agent aims to discover an optimal policy π^* that maximizes the expected cumulative reward over time through environment interaction and policy refinement based on received rewards [115].

B. Motivations and Rationales for Employing DRL

This work focuses on minimizing the peak AoI in an LBD-powered multi-UAV data collection system in IoT networks by optimizing UAV flight trajectories and charging strategies. The formulated problem exhibits the following key attributes:

- *Real-time Decision-making:* In the formulated problem, UAVs need to continuously adjust their flight paths and charging strategies in response to fluctuating AoI of IoTs, varying energy levels, and positions of neighboring UAVs, introducing considerable uncertainty into the decision-making process.
- *Long-term Optimization Objectives:* In the formulated problem, decisions made by UAVs at any given time step influence subsequent AoI and energy levels, requiring an approach that prioritizes long-term optimization goals rather than merely immediate benefits.
- *Computational Complexity:* As demonstrated in Section III, the challenge is the NP-hard nature, which extends beyond the traveling salesman problem with added complexities including multiple UAVs, energy constraints, charging strategies, and peak AoI minimization.

Conventional optimization methods struggle to deal with the above challenges of the formulated optimization problem. Firstly, the problem is with NP-hard nature and complex non-linear relationships between objectives and decision variables, which may render conventional optimization techniques ineffective (e.g., exhaustive approach or convex optimization [116], [117]). Secondly, the formulated problem requires sequential decision-making in a dynamic environment and trade-off among optimization objectives, which leads to poor performance of conventional algorithms such as evolutionary algorithms that rely on accurate prior knowledge [118]. Finally, the vast solution space generated by trajectory planning and charging decisions of multiple UAVs, combined with intricate energy consumption models, makes developing an efficient online algorithm impractical [119].

In this case, DRL provides important advantages for these optimization challenges, particularly in the dynamics of the considered scenario. First, DRL learns from environmental interactions through trial and error, thereby enabling adaptation to changing conditions. This powerful adaptation capability of DRL makes it especially effective in the considered scenarios where IoT data generation patterns, UAV energy states, and AoI are in constant flux. Moreover, the capacity of DRL to optimize for future rewards enables it to effectively balance conflicting objectives (i.e., peak AoI minimization and energy limitations) by considering long-term outcomes rather than

solely immediate gains. Consequently, the robust generalization capabilities of DRL and its proficiency in learning under uncertainty make it well-suited for the considered environment.

C. Necessary Principles of MADRL

MADRL enables multiple agents to learn optimal strategies through environment interaction and inter-agent communication, maximizing their long-term rewards without complete a priori knowledge. Different from the conventional DRL focuses on single-agent scenarios, MADRL extends this framework to multi-agent systems where decisions of agents mutually influence their rewards. Thus, POMDP is often adopted in multi-agent settings due to the limited observability of the agent.

A POMDP is formally defined as a tuple $(\mathcal{S}, \{\mathcal{A}_i\}_{i=1}^N, \mathcal{P}, \{\mathcal{R}_i\}_{i=1}^N, \{\mathcal{O}_i\}_{i=1}^N)$, where \mathcal{S} represents the state space, \mathcal{A}_i denotes action space of agent i , \mathcal{P} indicates the state transition probability function, \mathcal{R}_i defines immediate reward function of agent i , and \mathcal{O}_i describes observation space of agent i .

In MADRL, each agent i aims to find an optimal policy π_i^* that maximizes $\mathbb{E}_{\pi_i}(G_{i,t})$, where $\mathbb{E}_{\pi_i}(\cdot)$ represents the expected value under policy π_i . To determine π_i^* , a Q-value function $Q_i(s, a_1, a_2, \dots, a_N)$ is introduced, estimating expected discounted cumulative reward of agent i when actions (a_1, a_2, \dots, a_N) are executed at state s . In the following, we reformulate the optimization problem as a POMDP.

D. POMDP Formulation

Based on the MADRL principles mentioned above, we now formulate our optimization problem as a POMDP by defining its key components to address the UAV trajectory optimization and data collection challenges.

- **State space \mathcal{S} :** To fully capture the system dynamics, the state space must include all essential information for decision-making. Thus, we design the state space to consist of UAV positions for trajectory tracking, UAV energy levels for charging strategy, and AoI of IoTs and information status for data freshness monitoring. At time slot t , the system state is defined as

$$\mathcal{S}_t = \{ (x_j(t), y_j(t)), E_j(t), a_i(t), s_i(t) \mid i = 1, 2, \dots, N_S, j = 1, 2, \dots, N_U \}, \quad (6)$$

where $(x_j(t), y_j(t))$ represents the horizontal coordinate of UAV $u_j \in \mathcal{U}$.

- **Observation space \mathcal{O} :** In practical scenarios, each UAV has limited sensing capabilities and can only obtain local information. Different from the global state space, we design the observation space to reflect this partial observability. Each UAV observes its own position, energy level, and only the information status and AoI of accessible or nearby IoTs. Thus, the observation of UAV u_j at time slot t is defined as

$$\mathcal{O}_t^j = \{ (x_j(t), y_j(t)), E_j(t), a_i(t), s_i(t) \mid$$

$$i = 1, 2, \dots, N_S^j \}. \quad (7)$$

As such, the observation space for all UAVs at time slot t is denoted as

$$\mathcal{O}_t = \{ \mathcal{O}_t^j \mid j = 1, 2, \dots, N_U \}. \quad (8)$$

- **Action space \mathcal{A} :** To ensure practical implementation and reduce computational complexity, the horizontal movements of UAVs are limited to eight standard directions (i.e., north, northeast, east, etc.). This direction-based action design not only corresponds to realistic UAV control strategies, but also effectively captures the temporal dynamics of UAV motion. Furthermore, by constraining the action space to directional outputs rather than absolute positional coordinates, the bounded nature of the directional space enhances algorithmic convergence properties [120]–[123]. The action of UAV u_j at time slot t is defined as $\mathcal{A}_t^j = \{ a_j^x(t), a_j^y(t) \}$, and the action space for all UAVs is denoted as

$$\mathcal{A}_t = \{ \mathcal{A}_t^j \mid j = 1, 2, \dots, N_U \}. \quad (9)$$

- **Immediate reward \mathcal{R} :** To align with our optimization objective, the reward function is carefully designed to balance multiple goals. The immediate reward function r_t considers both the peak AoI of the network and a penalty term $r_p(t)$ defined as follows:

$$r_p(t) = \begin{cases} -d_c \cdot r_{pen1}, & \text{if } E_j(t) \leq E_\phi, \\ -d_c \cdot r_{pen2}, & \text{if } E_j(t) = E, \\ r_0, & \text{if } E_\phi < E_j(t) < E, \end{cases} \quad (10)$$

where d_c represents the distance from UAVs to the charging area boundary, E_ϕ denotes the UAV energy threshold for immediate charging, and E indicates the full UAV energy level. This penalty design encourages efficient energy management by penalizing both low energy states (risking operation interruption) and full energy states (indicating inefficient charging).

Following this, the comprehensive reward function is then designed as follows:

$$r_t = \alpha \cdot r_a(t) + \beta \cdot r_p(t) + \gamma \cdot r_s(t), \quad (11)$$

where $r_a(t)$ represents the reward associated with the AoI in the network at time slot t , $r_s(t)$ denotes the reward for successful information collection from IoTs at time slot t , α , β , and γ are weight factors for different parts of the reward function, respectively. As can be seen, this design directly relates to our optimization goal by penalizing high AoI while encouraging timely data collection and proper energy management.

Through this POMDP formulation, we can transform our optimization problem into a format suitable for MADRL algorithms. In the following, we introduce and analyze the standard MAPPO algorithm.

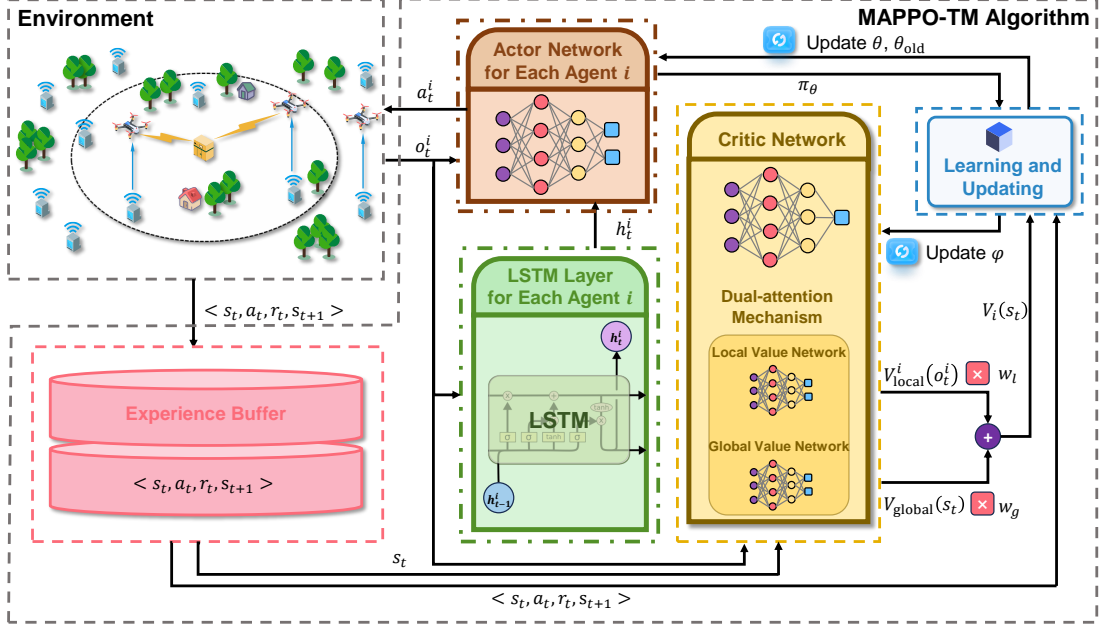


Fig. 2. The framework of the proposed MAPPO-TM algorithm.

E. Standard MAPPO Algorithm

In this work, we adopt MAPPO as our solution framework, which extends PPO to handle multi-agent scenarios effectively. The algorithm generates independent policy networks π_{θ_i} for each agent i , enabling decision-making based on their local observations. The optimization objective of MAPPO is given by

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T r_t \right], \quad (12)$$

where τ denotes the joint trajectories under policy π_{θ} , and r_t represents the reward at time step t .

To address the scalability issue in MADRL, MAPPO employs centralized training with decentralized execution (CTDE). During training, a centralized critic evaluates the global state by minimizing the temporal difference error:

$$J_{\text{critic}}(V) = \mathbb{E}_{\tau} \left[\sum_{t=0}^T \frac{1}{2} (r_t + \gamma V(s_{t+1}) - V(s_t))^2 \right], \quad (13)$$

where γ is the discount factor and $V(s_t)$ represents the global value function at time step t .

Following the PPO design, MAPPO adopts a clipped surrogate objective for stable policy updates as follows:

$$L_i^{\text{clip}}(\theta) = \mathbb{E}_t \left[\min(\rho_t^i(\theta) A_t^i, \text{clip}(\rho_t^i(\theta), 1 - \epsilon, 1 + \epsilon) A_t^i) \right], \quad (14)$$

where $\rho_t^i(\theta)$ denotes the probability ratio, A_t^i represents the advantage function for agent i , and ϵ controls the update step size.

However, the standard MAPPO algorithm faces two main challenges in our scenario. First, the specially designed reward function that balances AoI, energy management, and data collection may lead to unstable convergence during training.

Second, the complex coordination among multiple UAVs to achieve efficient coverage while avoiding energy depletion requires more sophisticated policy updates. Thus, we propose several enhancements to the algorithm in the following section to address these challenges.

F. MAPPO-TM Algorithm

Based on the standard MAPPO, we propose MAPPO-TM to address the specific challenges in our UAV-enabled data collection scenario.

The algorithm introduces two key enhancements, which are temporal dependency learning for capturing UAV historical trajectories and energy states, and multi-agent coordination mechanism for balancing individual and global objectives. These improvements effectively address the convergence and coordination issues faced by standard MAPPO in our POMDP formulation, and are detailed as follows.

1) *Temporal Dependency Learning*: In our POMDP formulation, the UAVs need to make decisions based on both current observations and historical information (e.g., previous trajectories and energy consumption patterns). However, standard MAPPO with feed-forward networks can only process current states, leading to suboptimal decisions in temporally dependent scenarios. To address this limitation, we incorporate long short-term memory (LSTM) networks into the actor network, so that allowing each UAV to maintain and utilize historical information effectively.

Specifically, the LSTM-enhanced actor network for agent i is designed as follows:

$$h_t^i = \text{LSTM}(o_t^i, h_{t-1}^i; \theta_{\text{LSTM}}^i), \quad (15)$$

$$\pi_{\theta_i}(a_t^i | o_t^i) = f_{\text{actor}}(h_t^i; \theta_{\text{actor}}^i), \quad (16)$$

where h_t^i represents the hidden state that captures temporal dependencies, o_t^i is the current observation, and f_{actor} denotes the policy head that outputs action probabilities.

By using LSTM-enhanced actor networks, we capture complex temporal dependencies across UAV trajectories, energy states, and data freshness metrics. This allows our agents to maintain critical historical context for decision optimization, unlike standard MAPPO which is limited to processing only current states through feed-forward networks. This enhancement enables each UAV to better plan its trajectory and energy usage by considering historical patterns.

2) *Multi-agent Coordination Mechanism*: Standard MAPPO faces challenges in balancing individual UAV objectives (energy management) with global goals (network AoI minimization) in our scenario. As depicted in Fig. 2, to enhance coordination, we introduce a dual-attention mechanism that allows each UAV to weigh both local and global information during decision-making. The coordination-enhanced value function is given by

$$V_i(s_t) = w_l V_{\text{local}}^i(o_t^i) + w_g V_{\text{global}}(s_t), \quad (17)$$

where w_l and w_g are learnable weights that balance local and global objectives, V_{local}^i evaluates individual performance, and V_{global} assesses the overall system state.

By implementing a dual-attention coordination mechanism, we enable each UAV to dynamically weigh local and global information during decision-making. This sophisticated coordination balances individual UAV objectives with system-wide performance metrics, particularly important when multiple UAVs must collectively optimize data freshness while managing individual energy constraints. This mechanism helps UAVs maintain energy efficiency while contributing to the global AoI optimization goal.

Through these enhancements, MAPPO-TM significantly improves upon standard MAPPO by enabling temporal-aware decision-making and better multi-agent coordination, leading to more efficient UAV trajectory planning and data collection strategies.

3) *Main Steps of MAPPO-TM Algorithm*: Following the POMDP formulation and algorithm design, we now present the main steps of MAPPO-TM, as shown in Algorithm 1. The algorithm begins by initializing the LSTM-based actor networks with parameters θ for decentralized execution, a centralized critic network with parameters ϕ for global state evaluation, and an old actor network with parameters θ_{old} for stable policy updates. For each episode, after environment reset and initial state s_t acquisition, each UAV i obtains its local observation o_t^i and generates action a_t^i through its LSTM-enhanced actor network, capturing temporal dependencies in the decision-making process. The environment then executes the joint action set $\mathbf{a}_t = \{a_t^1, a_t^2, \dots, a_t^N\}$, providing reward r_t and next state s_{t+1} . These interaction experiences, including states, actions, rewards, and next states, are stored in an experience buffer for subsequent learning.

The learning process consists of two key components, which are critic network updates and actor network updates. The centralized critic network is updated by minimizing the value function loss, which evaluates global state values while

Algorithm 1: MAPPO-TM

```

1 Initialize the parameters  $\theta$  of the LSTM-based actor
  networks and the parameters  $\phi$  of the centralized
  critic network ;      // Initialize temporal
  learning networks.
2 Initialize the parameters of the old actor networks
   $\theta_{\text{old}} \leftarrow \theta$  ;      // For stable updates.
3 for  $Episode = 1, \dots, N_{\text{eps}}$  do
4   Reset the environment and initialize the state  $s_t$  ;
    // Start new episode.
5   for  $Time\ slot\ t = 1, \dots, T$  do
6     for  $Each\ agent\ i$  do      // Decentralized
      execution phase.
7       Obtain the current state  $s_t^i$  ;      // UAV
        local observation.
8       Pass  $s_t^i$  through the LSTM-based actor to
        generate action  $a_t^i$  ;
        // Temporal-aware decisions.
9     Execute joint actions  $\mathbf{a}_t = \{a_t^1, \dots, a_t^N\}$  ;
      // UAVs trajectory execution.
10    Receive reward  $r_t$  and next state  $s_{t+1}$  ;
      // AoI and energy-based reward.
11    Store transitions  $(s_t, \mathbf{a}_t, r_t, s_{t+1})$  into the
      experience buffer;
12    Update the centralized critic using the value
      function loss ;      // With
      dual-attention mechanism.
13    Update the actor networks using the clipped
      surrogate objective;
14    Update  $\theta_{\text{old}} \leftarrow \theta$  periodically ;
      // Stabilize training.

```

considering both individual UAV performance and network-wide AoI optimization. Concurrently, the actor networks are updated using the clipped surrogate objective enhanced with our coordination mechanism, ensuring stable policy improvements while maintaining effective multi-agent cooperation. To further stabilize the training process, the old actor network parameters θ_{old} are periodically synchronized with the current parameters θ , thereby maintaining consistent clipping ratios in the surrogate objective. This iterative process continues for N_{eps} episodes, thus progressively improving both individual UAV policies and overall system performance.

By building upon the inherent stability with its clipped surrogate objective of PPO, we ensure controlled policy updates crucial for convergence in our complex environment. This stability foundation addresses the high dimensionality and stochastic elements of our problem domain, while our targeted enhancements overcome standard PPO limitations.

G. Complexity Analysis

The computational complexity of the MAPPO-TM algorithm can be analyzed in both the training and execution phases.

TABLE III
ALGORITHMS COMPLEXITY COMPARISON.

Algorithms	Time Complexity	Space Complexity
MAPPO-TM	$\mathcal{O}(N_{\text{eps}}TN \theta_a + N_{\text{eps}}T \theta_c)$	$\mathcal{O}(N \theta_a + \theta_c + D(s + a + 1))$
MAPPO	$\mathcal{O}(N_{\text{eps}}TN \theta'_a + N_{\text{eps}}T \theta'_c)$	$\mathcal{O}(N \theta'_a + \theta'_c + D(s + a + 1))$
MATD3	$\mathcal{O}(N_{\text{eps}}TN \theta''_a + N_{\text{eps}}TN2 \theta''_c)$	$\mathcal{O}(N \theta''_a + N2 \theta''_c + D(s + a + 1))$
MADDPG	$\mathcal{O}(N_{\text{eps}}TN \theta'''_a + N_{\text{eps}}TN \theta'''_c)$	$\mathcal{O}(2N \theta'''_a + 2N \theta'''_c + D(s + a + 1))$

- **For the training phase**, the time complexity consists of several key components. First, the network initialization requires $\mathcal{O}(N|\theta_a| + |\theta_c|)$ computations for both actor networks and centralized critic network, where $|\theta_a|$ and $|\theta_c|$ represent their respective parameter sizes. Second, the action sampling process through LSTM-based actor networks takes $\mathcal{O}(N_{\text{eps}}TN|\theta_a|)$ operations, where N_{eps} denotes the training episodes, T represents time slots per episode, and N is the number of UAVs. Third, the experience storage requires $\mathcal{O}(N_{\text{eps}}TNV)$ operations for maintaining the replay buffer, where V denotes the size of state-action-reward tuples. Finally, the network update process needs $\mathcal{O}(N_{\text{eps}}T|\theta_c| + N_{\text{eps}}TN|\theta_a|)$ computations for both critic and actor networks. Thus, the overall time complexity during training is $\mathcal{O}(N_{\text{eps}}TN|\theta_a| + N_{\text{eps}}T|\theta_c|)$.

The space complexity during training mainly comes from two aspects, which are network parameters storage and replay buffer maintenance. The former requires $\mathcal{O}(N|\theta_a| + |\theta_c|)$ space for storing all network parameters, while the latter needs $\mathcal{O}(D(|s| + |a| + 1))$ space for maintaining the replay buffer with size D , where $|s|$ and $|a|$ represent the dimensions of state and action spaces, respectively. Therefore, the total space complexity in training phase is $\mathcal{O}(N|\theta_a| + |\theta_c| + D(|s| + |a| + 1))$.

- **For the execution phase**, since only the actor networks are utilized for action sampling, the time complexity reduces to $\mathcal{O}(TN|\theta_a|)$ and the space complexity becomes $\mathcal{O}(N|\theta_a|)$. This significant reduction in computational complexity makes MAPPO-TM suitable for real-world UAV deployments.

In addition, the complexity comparison between the proposed MAPPO-TM algorithm and other conventional MADRL algorithms is shown in Table III. As can be seen, MAPPO-TM offers a balanced computational complexity compared to other MADRL methods. The efficiency-performance balance makes MAPPO-TM particularly well-suited for the UAV trajectory optimization and data collection challenges addressed in this work.

V. SIMULATION RESULTS

In this section, we present simulation results and analyses. We first introduce the simulation setup and benchmarks, and then provide simulation results.

A. Simulation Setups

We consider a typical field environment with multiple IoTs and LBDs. The UAVs fly in the airspace above the square area. The scene is divided into a circular charging area at the center and the rest of the non-charging area. The coordinate system takes the center of the square area, i.e., the coordinates of the LBD, as the origin. The initial coordinates of the IoTs are fixed, while the initial coordinates of the four UAVs are $(1, 0, H)$, $(-1, 0, H)$, $(0, 1, H)$, and $(0, -1, H)$, respectively. Furthermore, the initial power of the UAVs is 60% of the full power. Other parameters are shown in Table IV.

For comparison, we utilize MAPPO [124], multi-agent twin delayed deep deterministic policy gradient (MATD3) [125], and MADDPG [126] as benchmark methods.

- **MAPPO**: MAPPO extends the PPO framework to multi-agent environments. It employs CTDE, allowing each agent to learn its policy while leveraging global information during training. MAPPO maintains the stability and sample efficiency of PPO, making it suitable for complex multi-agent tasks with continuous action spaces. Additionally, it incorporates mechanisms to handle the non-stationarity arising from multiple learning agents, ensuring robust and scalable policy learning.
- **MATD3**: MATD3 adapts the twin delayed deep deterministic policy gradient (TD3) algorithm for multi-agent settings. It utilizes multiple critic networks to mitigate overestimation biases and employs a delayed update strategy to stabilize training across agents. MATD3 leverages decentralized policies while maintaining centralized critics, enabling effective coordination among agents in environments with high noise and uncertainty. This approach enhances the accuracy of value estimates and promotes cooperative behavior in multi-agent scenarios.
- **MADDPG**: MADDPG extends the DDPG algorithm to multi-agent systems. It employs CTDE, where each agent has its own actor network and its own centralized critic network that considers the observations and actions of all agents. This framework helps address the non-stationarity problem in multi-agent environments by providing each agent with additional information during training. MADDPG is effective in continuous action spaces and facilitates coordinated strategies among multiple agents, improving overall performance in complex tasks.

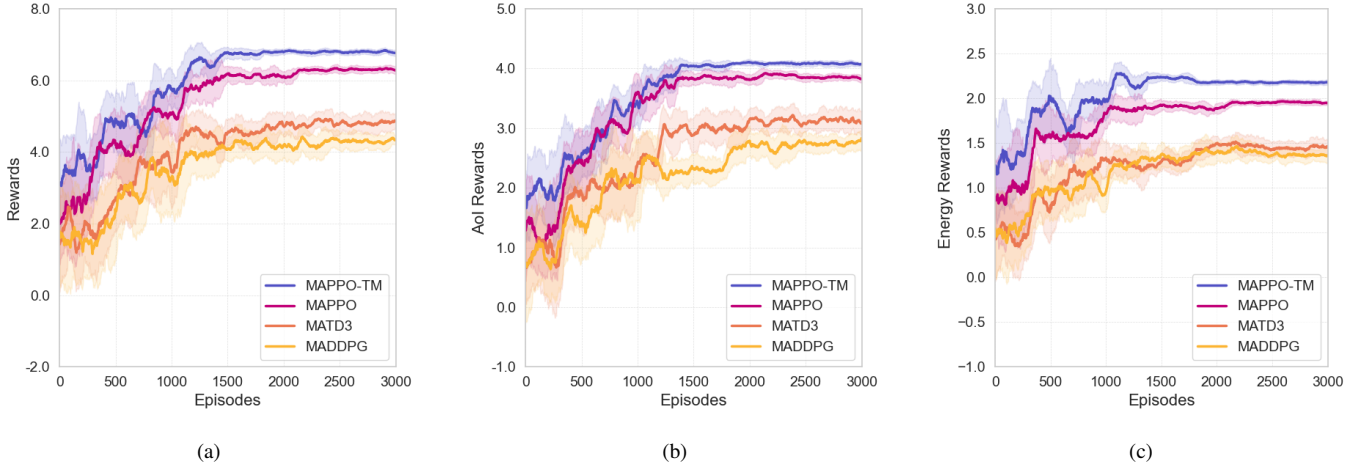


Fig. 3. Training results. (a) Cumulative rewards training curve. (b) AoI rewards training curve. (c) Energy rewards training curve.

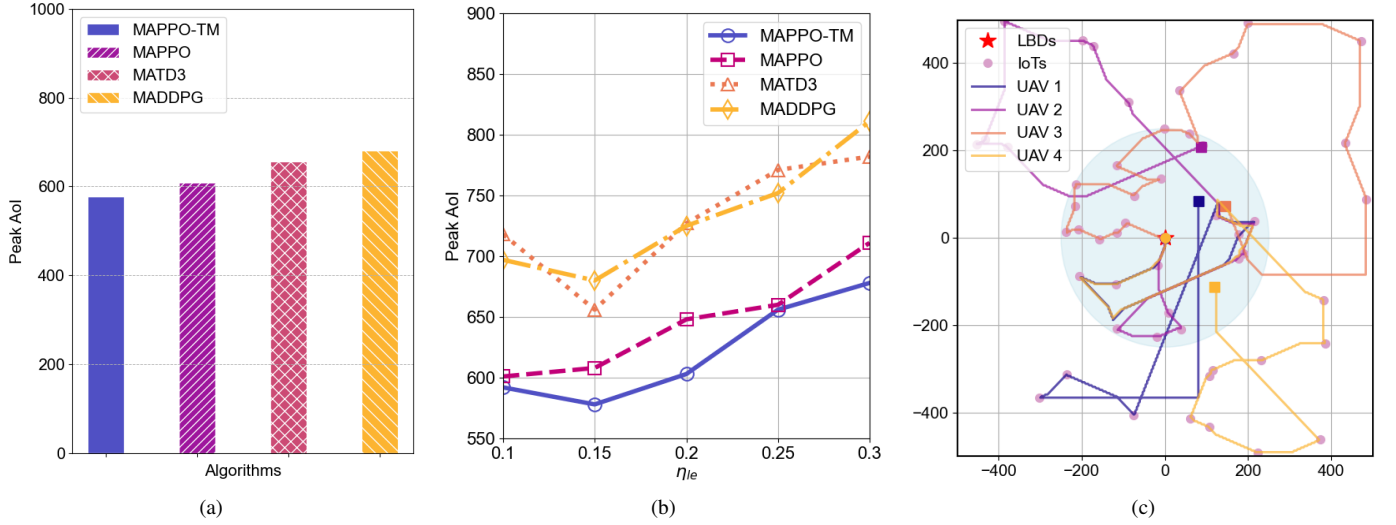


Fig. 4. (a) Peak AoI comparison across algorithms. (b) Peak AoI changing across different values of η_{le} . (c) Flight Paths of UAVs driven by MAPPO-TM

B. Simulation Results

In this subsection, we evaluate the performance of the proposed MAPPO-TM framework under various UAV-enabled data collection scenarios.

1) *Reward Performance Analysis*: Fig. 3(a) presents a comparative analysis of the four algorithms in terms of cumulative rewards. The results demonstrate that during the early training stages, the proposed MAPPO-TM algorithm maintains relatively stable performance at a higher reward level, while the conventional MAPPO algorithm performs marginally below MAPPO-TM. In contrast, MATD3 and MADDPG exhibit considerable fluctuations at lower reward levels, with MADDPG occasionally experiencing negative rewards. As training progresses, MAPPO-TM significantly outperforms the other algorithms. This superior performance can be attributed to the LSTM-based temporal dependency learning mechanism, which enables each UAV to effectively maintain and utilize historical information, thereby enhancing reward stability during initial training phases. Furthermore, the dual-attention

coordination mechanism guides UAVs to balance local and global objectives, thus optimizing network-wide performance and maximizing cumulative rewards.

Fig. 3(b) depicts the AoI rewards over 3000 episodes for MAPPO-TM, MAPPO, MATD3 and MADDPG. MAPPO-TM significantly outperforms the baselines, stabilizing at an AoI reward of approximately 4.1 after 1500 episodes, while MAPPO plateaus at around 3.9, MATD3 at 3.2, and MADDPG at 2.6. The shaded regions around each line indicate variance, with MAPPO-TM showing the smallest variance, reflecting stable performance, whereas MAPPO, MATD3, and MADDPG exhibit larger fluctuations, particularly in the early stages, with MADDPG displaying the highest variability throughout. The superior AoI performance of MAPPO-TM is primarily due to its multi-agent coordination mechanism, which uses a dual-attention approach to balance individual UAV objectives with the global goal of AoI minimization, thereby enabling more effective data collection strategies. Additionally, the LSTM-based temporal dependency learning allows MAPPO-TM to leverage historical trajectories, addressing the limita-

TABLE IV
SIMULATION PARAMETERS.

Parameter	Value	Parameter	Value
A	0.1256 m ²	R_c	250 m
d_0	0.5009	t_d	1 s
E_u	30000 J	v	5 m/s
H	80 m	v_0	5.0463 m/s
N_D	10	v_{tip}	80 m/s
N_S	50	δ	10 ⁻⁶
N_U	4	η_{le}	0.15
P_L	1000 W	ρ	1.225 kg/m ³
P_α	14.7517 W	$\frac{\beta_0}{\sigma^2}$	80 dB
P_β	41.5409 W	ω	0.1248

tions of the feed-forward networks in standard MAPPO, which struggle with long-term planning, and the lack of coordination in MATD3 and MADDPG, thereby leading to their lower AoI rewards.

Fig. 3(c) illustrates the energy rewards over 3000 episodes for the same four algorithms. MAPPO-TM again achieves the highest energy rewards, stabilizing at around 2.2 after 1800 episodes, while MAPPO plateaus at 2.0, MATD3 at 1.4, and MADDPG at 1.3. The variance, indicated by the shaded regions, is smallest for MAPPO-TM, suggesting consistent energy efficiency, whereas MAPPO shows moderate fluctuations, and MATD3 and MADDPG exhibit significant variability, with MADDPG performing the least reliably. The enhanced energy efficiency of MAPPO-TM can be attributed to its LSTM-enhanced actor network, which captures historical energy consumption patterns, allowing UAVs to make informed decisions that optimize energy usage over time, different from the baselines that lack temporal awareness. Furthermore, the coordination mechanism of MAPPO-TM ensures that UAVs balance energy management with global objectives, thus preventing the energy depletion issues seen in MATD3 and MADDPG, which struggle to coordinate effectively, and in MAPPO, which cannot account for long-term energy dependencies, thereby resulting in their lower energy rewards.

2) *Peak AoI Performance Evaluation:* To further validate the effectiveness of MAPPO-TM, we analyze peak AoI performance across the different algorithms. Fig. 4(a) illustrates the final peak AoI values, with MAPPO-TM achieving significantly lower AoI compared to MAPPO and substantially outperforming other algorithms, while MATD3 and MADDPG demonstrate considerably higher peak AoI values respectively. This superior performance of MAPPO-TM can be attributed to two key factors. First, the LSTM network effectively captures temporal dependencies in UAV trajectories and energy states, enabling more informed decision-making and reducing the probability of suboptimal actions, consequently lowering peak AoI. Second, the dual-attention mechanism successfully balances individual UAV objectives with global

network performance, thus guiding UAVs toward trajectories that simultaneously optimize energy efficiency and network-wide AoI.

Additionally, Fig. 4(b) illustrates peak AoI variations corresponding to different values of η_{le} , where η_{le} represents the laser-to-electricity conversion efficiency in Eq. (2). Notably, almost all algorithms achieve optimal performance at $\eta_{le} = 0.15$. When η_{le} falls below 0.15, the reduced laser charging power extends UAV charging time, limiting their ability to collect IoT information outside the charging zone, thereby increasing AoI. Conversely, higher η_{le} values result in shorter charging durations, preventing sufficient data collection within the charging zone, thus also increasing AoI. MAPPO-TM consistently outperforms other algorithms across different η_{le} values, achieving a 5-10% lower peak AoI. This superior performance can be attributed to its LSTM mechanism, which effectively captures temporal dependencies in the UAV charging and information collection processes, thereby enabling more informed decision-making and better system optimization.

3) *UAV Flight Trajectory Analysis:* Fig. 4(c) depicts the flight trajectories of multiple UAVs in the system. The UAVs start near the LBDs with partially charged batteries and initially collect data of IoTs in the charging area. Then UAVs fly out to perform tasks once fully charged, and return to recharge when battery levels drop below a threshold, thereby demonstrating efficient task division with minimal overlap. This optimized behavior is driven by LSTM-based temporal dependency learning of MAPPO-TM, which captures historical trajectories and energy patterns for better long-term planning, and its multi-agent coordination mechanism, which uses a dual-attention approach to balance individual energy management with global AoI minimization. The stable policy updates of the algorithm, facilitated by a clipped surrogate objective and periodic parameter synchronization, further ensure consistent trajectory decisions, thus preventing erratic movements. As a result, MAPPO-TM enhances mission endurance, reduces AoI, and improves overall efficiency.

In summary, the proposed MAPPO-TM algorithm consistently outperforms alternative algorithms in terms of cumulative rewards and peak AoI optimization in multi-UAV data collection scenarios. Moreover, the LSTM-based temporal dependency learning enhances decision-making stability, while the dual-attention coordination mechanism significantly improves multi-agent cooperation, thereby resulting in superior overall system performance.

VI. DISCUSSION

In this section, we present some discussions related to the system and algorithm.

A. Scalability Analysis

The considered LBD-powered multi-UAV data collection system is designed with adaptability and flexibility at its core, enabling it to function effectively across various IoT deployments. The following key aspects demonstrate this adaptability:

- *The considered system is fundamentally agnostic to the specific hardware characteristics of IoTs, focusing instead on their functional capabilities for data generation and communication.* This design principle allows our approach to work with heterogeneous IoT deployments spanning multiple application domains. While IoTs may vary significantly in their processing power, memory capacity, and energy resources, our system primarily requires them to maintain basic communication capabilities and data buffering functionality. The UAV trajectory planning and laser charging scheduling optimized by MAPPO-TM operates independently of the internal complexities of individual IoTs, instead focusing on their spatial distribution, data generation patterns, and AoI characteristics.
- *The computational complexity of our MAPPO-TM algorithm has been carefully analyzed and optimized, as detailed in Section IV of our paper.* During the execution phase, the time complexity reduces to $\mathcal{O}(TN|\theta_a|)$ and the space complexity becomes $\mathcal{O}(N|\theta_a|)$, where T represents the number of time slots, N is the number of UAVs, and $|\theta_a|$ denotes the parameter size of the actor networks. This significant reduction in computational requirements during deployment enables our system to function efficiently even when computational resources are constrained. The LSTM-based temporal memory mechanism in our algorithm is particularly efficient, as it allows for selective retention of only the most relevant historical information, thereby minimizing memory requirements while maximizing decision quality.
- *The proposed MAPPO-TM framework employs a CTDE paradigm, which provides substantial deployment flexibility.* The centralized training can be performed on powerful computing infrastructure (e.g., cloud servers or edge computing nodes), while the trained models are deployed to individual UAVs for decentralized execution. This approach effectively separates the computationally intensive training process from the resource-constrained deployment environment. Our simulation results demonstrate that MAPPO-TM maintains stable performance across different laser-to-electricity conversion efficiency values (η_{le}), highlighting its robustness to variations in the physical characteristics of the deployment environment.

As such, the MAPPO-TM framework represents a flexible and adaptable solution for AoI optimization in laser-charged UAV-assisted IoT networks. The ability of the MAPPO-TM framework to function effectively across diverse IoT device types, computational environments, and network configurations makes it well-suited for practical deployment in various real-world scenarios.

B. Robustness Analysis

In practical deployments, UAV failures due to hardware malfunctions or environmental factors can significantly impact system performance. Our framework implements specific mechanisms to handle various failure scenarios at a granular

level. These mechanisms are designed to address failures of both IoT devices and UAVs during operation, ensuring system resilience through adaptive behavior rather than requiring complete system reconfiguration. The integration of these failure handling techniques directly into our MADRL framework allows for real-time response to unexpected events without compromising overall mission objectives. The details are as follows:

- *Faulty IoT Handling:* In our system, when a UAV detects that an IoT is unresponsive or malfunctioning (i.e., no data transmission is occurring despite being within communication range), the UAV updates its observation space \mathcal{O}_t^i to mark this IoT as faulty. This information is then shared with other UAVs during the centralized training phase. The MAPPO-TM algorithm is designed to adapt dynamically to such changes by adjusting the UAV trajectories to prioritize functioning IoTs, thereby maintaining efficient data collection despite IoT failures. Additionally, the temporal memory mechanism in our LSTM-based actor network enables the UAVs to remember which IoTs have previously been identified as faulty, preventing repeated unsuccessful visit attempts in subsequent time slots [127], [128].
- *UAV Failure During Operation:* Our multi-agent framework has built-in resilience against UAV failures. If a UAV becomes inoperable during a mission, the remaining UAVs can autonomously redistribute their responsibilities through the multi-agent coordination mechanism. Specifically, the dual-attention value function $V_i(s_t) = w_l V_{\text{local}}^i(o_t^i) + w_g V_{\text{global}}(s_t)$ allows UAVs to dynamically adjust their behavior based on global objectives when the system state changes due to a UAV failure. The system detects a UAV failure when it stops communicating its state updates, and the remaining UAVs then update their global value assessment to account for the reduced fleet capacity, resulting in an automatic reallocation of data collection responsibilities among the functioning UAVs [129], [130].
- *Energy-related UAV Failures:* Since energy management is a critical aspect of our system, we have specifically addressed UAV failures related to energy depletion. The penalty term $r_p(t)$ in our reward function (Eq. 10) strongly discourages UAVs from reaching critically low energy states. If a UAV energy level approaches a critical threshold despite these preventive measures, the system initiates a fail-safe protocol: the UAV immediately prioritizes returning to the nearest laser charging area while transmitting its current data payload to neighboring UAVs if possible. This approach ensures that even if a UAV cannot complete its mission due to energy constraints, the collected data is not lost, and the peak AoI performance is preserved as much as possible [131], [132].

Moreover, we can also incorporate some additional robustness measures into the considered system to further improve the ability to resist failures while minimizing the disruption to the data collection service. The specific implementation details of these robustness measures are as follows:

Firstly, we address the feasible retraining approaches that enable rapid algorithm convergence when UAV failures occur. Specifically, in the event of UAV malfunction, system administrators can adjust the UAV count parameters in the simulation environment and subsequently retrain the DRL algorithm. Given that this process can be executed in computationally robust settings, administrators may proactively develop multiple versions of the DRL algorithm trained with varying UAV quantities as contingency measures. During this retraining process, advanced techniques such as transfer learning [133] can significantly accelerate algorithm convergence. Furthermore, by utilizing edge computing infrastructure and implementing incremental learning methodologies [134], we can perform dynamic updates to the deployed neural networks, thereby enhancing UAV adaptability to fluctuations in operational fleet size.

Secondly, we highlight that redundancy-based fault-tolerance mechanisms offer substantial benefits for such emergency management. For example, maintaining a reserve fleet of standby UAVs allows for prompt deployment to replace malfunctioning units with minimal operational disruption. Such redundancy mechanisms facilitate real-time adaptation and preserve communication system performance without significant service interruptions.

These fault-handling mechanisms significantly enhance the robustness of our MAPPO-TM framework in real-world scenarios where hardware failures and energy management challenges are inevitable. The ability to dynamically adapt to such edge cases is a key advantage of our MADRL approach compared to traditional optimization methods, which typically require predefined contingency plans for each possible failure scenario.

VII. CONCLUSION

This paper has investigated an LBD-powered multi-UAV data collection system in IoT networks. By utilizing LBD to charge UAVs, we effectively improved the flight time of UAVs, thereby enabling the IoT data collection tasks to be completed in a more timely manner. Specifically, we formulated a joint optimization problem that aims to minimize the peak AoI of the IoT network while ensuring that the UAVs follow the charging strategy and do not run out of power, which is characterized by high real-time and dynamic complexity. To address this problem, we proposed a MAPPO-TM algorithm that integrates LSTM with the MAPPO algorithm, which incorporates time-dependent learning to capture the historical trajectory and energy state of the UAV. Furthermore, the algorithm implements a multi-agent coordination mechanism to balance individual and global objectives, thus improving stability and convergence speed. Simulation results validated the effectiveness of the proposed MAPPO-TM algorithm, which demonstrates superior performance compared to other baselines, particularly in enhancing learning stability and reducing peak AoI, and thus highlighting its potential for practical deployment in the LBD-powered multi-UAV data collection system in IoT networks.

To further enhance the robustness and applicability of the considered system and proposed algorithm in practical scenar-

ios, future works will focus on the following aspects. Firstly, future works include developing an adaptive fault-tolerance mechanism that can dynamically reconfigure the multi-UAV network topology when failures occur, incorporating rapid retraining methods with transfer learning to minimize disruption to data collection operations. Secondly, we plan to enhance the laser charging efficiency by investigating advanced beam tracking algorithms and optimizing the trade-off between charging time and data collection to further reduce peak AoI in larger-scale deployments. Finally, we will incorporate domain randomization techniques during training to improve the robustness of the MAPPO-TM algorithm against environmental uncertainties and extend its generalization capabilities to diverse operational conditions.

REFERENCES

- [1] X. Liu, Z. Liu, B. Lai, B. Peng, and T. S. Durrani, "Fair energy-efficient resource optimization for multi-UAV enabled Internet of Things," *IEEE Trans. Veh. Technol.*, vol. 72, no. 3, pp. 3962–3972, 2023.
- [2] X. Liu, B. Lai, B. Lin, and V. C. M. Leung, "Joint communication and trajectory optimization for multi-UAV enabled mobile Internet of Vehicles," *IEEE Trans. Intell. Transp. Syst.*, vol. 23, no. 9, pp. 15 354–15 366, 2022.
- [3] M. Lahmeri, M. A. Kishk, and M. Alouini, "Charging techniques for UAV-assisted data collection: Is laser power beaming the answer?" *IEEE Commun. Mag.*, vol. 60, no. 5, pp. 50–56, 2022.
- [4] F. H. Panahi and F. H. Panahi, "Reliable and energy-efficient UAV communications: A cost-aware perspective," *IEEE Trans. Mob. Comput.*, vol. 23, no. 5, pp. 4038–4049, 2024.
- [5] J. Liu, X. Wang, B. Bai, and H. Dai, "Age-optimal trajectory planning for UAV-assisted data collection," in *Proc. IEEE INFOCOM*, 2018, pp. 553–558.
- [6] F. Song, M. Deng, H. Xing, Y. Liu, F. Ye, and Z. Xiao, "Energy-efficient trajectory optimization with wireless charging in UAV-assisted MEC based on multi-objective reinforcement learning," *IEEE Trans. Mob. Comput.*, vol. 23, no. 12, pp. 10 867–10 884, 2024.
- [7] I. Benmad, E. Driouch, and M. Kardouchi, "Data collection in UAV-assisted wireless sensor networks powered by harvested energy," in *Proc. IEEE PIMRC*, 2021, pp. 1351–1356.
- [8] C. Zhang, J. Liu, L. Xie, and X. He, "Age-optimal data gathering and energy recharging of UAV in wireless sensor networks," in *Proc. ACM AISS*, 2021, pp. 78:1–78:6.
- [9] A. Ranjha and G. Kaddoum, "URLLC-enabled by laser powered UAV relay: A quasi-optimal design of resource allocation, trajectory planning and energy harvesting," *IEEE Trans. Veh. Technol.*, vol. 71, no. 1, pp. 753–765, 2022.
- [10] K. Ahmadi and W. A. Serdijn, "Advancements in laser and LED-based optical wireless power transfer for IoT applications: A comprehensive review," *IEEE Internet Things J.*, 2025.
- [11] K. Zhu, J. Yang, Y. Zhang, J. Nie, W. Y. B. Lim, H. Zhang, and Z. Xiong, "Aerial refueling: Scheduling wireless energy charging for UAV enabled data collection," *IEEE Trans. Green Commun. Netw.*, vol. 6, no. 3, pp. 1494–1510, 2022.
- [12] Y. Zeng, J. Xu, and R. Zhang, "Energy minimization for wireless communication with rotary-wing UAV," *IEEE Trans. Wirel. Commun.*, vol. 18, no. 4, pp. 2329–2345, 2019.
- [13] P. Wu, F. Xiao, H. Huang, and R. Wang, "Load balance and trajectory design in multi-UAV aided large-scale wireless rechargeable networks," *IEEE Trans. Veh. Technol.*, vol. 69, no. 11, pp. 13 756–13 767, 2020.
- [14] Z. Lu, X. Wang, Z. Wang, and Y. Pei, "Optimal scheduling and deep reinforcement learning for multimodal charging system via unmanned aerial vehicles," *IEEE Trans. Green Commun. Netw.*, 2024.
- [15] H. Hu, K. Xiong, G. Qu, Q. Ni, P. Fan, and K. B. Letaief, "AoI-minimal trajectory planning and data collection in UAV-assisted wireless powered IoT networks," *IEEE Internet Things J.*, vol. 8, no. 2, pp. 1211–1223, 2021.
- [16] O. A. Amodu, C. Jarray, R. A. R. Mahmood, H. Althumali, U. A. Bukar, R. Nordin, N. F. Abdullah, and N. C. Luong, "Deep reinforcement learning for AoI minimization in UAV-aided data collection for WSN and IoT applications: A survey," *IEEE Access*, vol. 12, pp. 108 000–108 040, 2024.

- [17] Y. Long, S. Zhao, S. Gong, B. Gu, D. Niyato, and X. Shen, "AoI-aware sensing scheduling and trajectory optimization for multi-UAV-assisted wireless backscatter networks," *IEEE Trans. Veh. Technol.*, vol. 73, no. 10, pp. 15 440–15 455, 2024.
- [18] D. Guo, K. Nakhleh, I. Hou, S. Kompella, and C. Kam, "AoI, timely-throughput, and beyond: A theory of second-order wireless network optimization," *IEEE/ACM Trans. Netw.*, vol. 32, no. 6, pp. 4707–4721, 2024.
- [19] Z. Shen, G. Zhou, H. Huang, C. Huang, Y. Wang, and F. Wang, "Convex optimization-based trajectory planning for quadrotors landing on aerial vehicle carriers," *IEEE Trans. Intell. Veh.*, vol. 9, no. 1, pp. 138–150, 2024.
- [20] R. Zhang, H. Du, Y. Liu, D. Niyato, J. Kang, S. Sun, X. Shen, and H. V. Poor, "Interactive AI with retrieval-augmented generation for next generation networking," *IEEE Netw.*, vol. 38, no. 6, pp. 414–424, 2024.
- [21] R. Zhang, H. Du, D. Niyato, J. Kang, Z. Xiong, A. Jamalipour, P. Zhang, and D. I. Kim, "Generative AI for space-air-ground integrated networks," *IEEE Wirel. Commun.*, vol. 31, no. 6, pp. 10–20, 2024.
- [22] Q. Chen, Z. Guo, W. Meng, S. Han, C. Li, and T. Q. S. Quek, "A survey on resource management in joint communication and computing-embedded SAGIN," *IEEE Commun. Surv. Tutorials*, 2024.
- [23] M. Yi, X. Wang, J. Liu, Y. Zhang, and B. Bai, "Deep reinforcement learning for fresh data collection in UAV-assisted IoT networks," in *Proc. IEEE INFOCOM*, 2020, pp. 716–721.
- [24] K. Li, W. Ni, E. Tovar, and A. Jamalipour, "On-board deep Q-network for UAV-assisted online power transfer and data collection," *IEEE Trans. Veh. Technol.*, vol. 68, no. 12, pp. 12 215–12 226, 2019.
- [25] J. Miao, S. Bai, S. Mumtaz, Q. Zhang, and J. Mu, "Utility-oriented optimization for video streaming in UAV-aided MEC network: A DRL approach," *IEEE Trans. Green Commun. Netw.*, vol. 8, no. 2, pp. 878–889, 2024.
- [26] B. I. Ghomri, M. Y. Bendimerad, and F. T. Bendimerad, "DRL-driven optimization for energy efficiency and fairness in NOMA-UAV networks," *IEEE Commun. Lett.*, vol. 28, no. 5, pp. 1048–1052, 2024.
- [27] T. Li, K. Zhu, N. C. Luong, D. Niyato, Q. Wu, Y. Zhang, and B. Chen, "Applications of multi-agent reinforcement learning in future Internet: A comprehensive survey," *IEEE Commun. Surv. Tutorials*, vol. 24, no. 2, pp. 1240–1279, 2022.
- [28] D. Alsadie, "Efficient task offloading strategy for energy-constrained edge computing environments: A hybrid optimization approach," *IEEE Access*, vol. 12, pp. 85 089–85 102, 2024.
- [29] Y.-S. Liao, Y.-W. P. Hong, and J.-P. Sheu, "Laser-powered UAV trajectory and charging optimization for sustainable data-gathering in the Internet of Things," *IEEE Trans. Mob. Comput.*, 2024.
- [30] W. Liu, S. Zhang, and N. Ansari, "Joint laser charging and DBS placement for drone-assisted edge computing," *IEEE Trans. Veh. Technol.*, vol. 71, no. 1, pp. 780–789, 2022.
- [31] M. Lahmeri, M. A. Kishk, and M. Alouini, "Laser-powered UAVs for wireless communication coverage: A large-scale deployment strategy," *IEEE Trans. Wirel. Commun.*, vol. 22, no. 1, pp. 518–533, 2023.
- [32] D. H. Nguyen, "Dynamic optical wireless power transfer for electric vehicles," *IEEE Access*, vol. 11, pp. 2787–2795, 2023.
- [33] L. Zhang, Y. Wang, M. Min, C. Guo, V. Sharma, and Z. Han, "Privacy-aware laser wireless power transfer for aerial multi-access edge computing: A Colonel Blotto game approach," *IEEE Internet Things J.*, vol. 10, no. 7, pp. 5923–5939, 2023.
- [34] Z. Cheng, Z. Gao, M. Liwang, L. Huang, X. Du, and M. Guizani, "Intelligent task offloading and energy allocation in the UAV-aided mobile edge-cloud continuum," *IEEE Netw.*, vol. 35, no. 5, pp. 42–49, 2021.
- [35] T. Du, X. Gui, X. Teng, K. Zhang, and D. Ren, "Dynamic trajectory design and bandwidth adjustment for energy-efficient UAV-assisted relaying with deep reinforcement learning in MEC IoT system," *IEEE Internet Things J.*, vol. 11, no. 23, pp. 37 463–37 479, 2024.
- [36] D. Li, S. Xu, C. Zhao, Y. Wang, R. Xu, and B. Ai, "Data collection in laser-powered UAV-assisted IoT networks: Phased scheme design based on improved clustering algorithm," *IEEE Trans. Green Commun. Netw.*, vol. 8, no. 1, pp. 482–497, 2024.
- [37] M. S. Bashir and M. Alouini, "Energy optimization of a laser-powered hovering-UAV relay in optical wireless backhaul," *IEEE Trans. Wirel. Commun.*, vol. 22, no. 5, pp. 3216–3230, 2023.
- [38] W. Liu, L. Zhang, and N. Ansari, "Laser charging enabled DBS placement for downlink communications," *IEEE Trans. Netw. Sci. Eng.*, vol. 8, no. 4, pp. 3009–3018, 2021.
- [39] Y. Park, D. Kim, and J. Lee, "Joint trajectory and charging power optimization for laser-charged UAV relaying networks," in *Proc. IEEE ICTC*, 2022, pp. 224–229.
- [40] X. Ma, X. Liu, and N. Ansari, "Green laser-powered UAV far-field wireless charging and data backhauling for a large-scale sensor network," *IEEE Internet Things J.*, vol. 11, no. 19, pp. 31 932–31 946, 2024.
- [41] A. M. Abdelhady, A. Çelik, C. Diaz-Vilor, H. Jafarkhani, and A. M. Eltawil, "Operation optimization of laser-powered aerial data harvesting for passive IoT networks," in *Proc. IEEE WCNC*, 2024, pp. 1–6.
- [42] X. Zhang, Y. Zhao, H. You, K. Jian, and L. Liang, "Resource allocation strategy for wireless powered communication networks with UAV-assisted edge computing," in *Proc. IEEE VTC*, 2024, pp. 1–6.
- [43] L. Wang, Y. Li, Y. Chen, T. Li, and Z. Yin, "Air-ground coordinated MEC: Joint task, time allocation and trajectory design," *IEEE Trans. Veh. Technol.*, vol. 74, no. 3, pp. 4728–4743, 2025.
- [44] K. Singh, P. Raut, P. K. Sharma, and C. Li, "Laser-powered multi-UAV URLLC systems: Reliability and scheduling performance analysis," *IEEE Trans. Veh. Technol.*, vol. 72, no. 11, pp. 14 615–14 630, 2023.
- [45] S. K. Singh, K. Agrawal, K. Singh, A. Bansal, C. Li, and Z. Ding, "On the performance of laser-powered UAV-assisted SWIPT enabled multiuser communication network with hybrid NOMA," *IEEE Trans. Commun.*, vol. 70, no. 6, pp. 3912–3929, 2022.
- [46] X. Wang, J. Li, Z. Ning, Q. Song, L. Guo, and A. Jamalipour, "Wireless powered metaverse: Joint task scheduling and trajectory design for multi-devices and multi-UAVs," *IEEE J. Sel. Areas Commun.*, vol. 42, no. 3, pp. 552–569, 2024.
- [47] Y. Fu, H. Mei, K. Wang, and K. Yang, "Joint optimization of 3D trajectory and scheduling for solar-powered UAV systems," *IEEE Trans. Veh. Technol.*, vol. 70, no. 4, pp. 3972–3977, 2021.
- [48] A. Goel and S. De, "Backscatter communication based sensor data collection using laser powered UAV," in *Proc. IEEE ICC*, 2023, pp. 2896–2901.
- [49] M. L. Betalo, S. Leng, A. M. Seid, H. N. Abishu, A. Erbad, and X. Bai, "Dynamic charging and path planning for UAV-powered rechargeable WSNs using multi-agent deep reinforcement learning," *IEEE Trans. Autom. Sci. Eng.*, 2025.
- [50] G. Zhang, X. Wei, X. Tan, Z. Han, and G. Zhang, "AoI minimization based on deep reinforcement learning and matching game for IoT information collection in SAGIN," *IEEE Trans. Commun.*, 2025.
- [51] B. Zhu, E. Bedeer, H. H. Nguyen, R. Barton, and Z. Gao, "UAV trajectory planning for AoI-minimal data collection in UAV-aided IoT networks by transformer," *IEEE Trans. Wirel. Commun.*, vol. 22, no. 2, pp. 1343–1358, 2023.
- [52] Q. Dang, Q. Cui, Z. Gong, X. Zhang, X. Huang, and X. Tao, "AoI oriented UAV trajectory planning in wireless powered IoT networks," in *Proc. IEEE WCNC*, 2022, pp. 884–889.
- [53] X. Xiao, X. Wang, and W. Lin, "Joint AoI-aware UAVs trajectory planning and data collection in UAV-based IoT systems: A deep reinforcement learning approach," *IEEE Trans. Consumer Electron.*, vol. 70, no. 4, pp. 6484–6495, 2024.
- [54] Q. Wang, X. Liang, H. Zhang, and L. Ge, "AoI-aware energy efficiency resource allocation for integrated satellite-terrestrial IoT networks," *IEEE Trans. Green Commun. Netw.*, vol. 9, no. 1, pp. 125–139, 2025.
- [55] J. Cao, X. Zhu, S. Sun, E. Kurniawan, and A. Boonkajay, "Risk-aware and energy-efficient AoI optimization for multi-connectivity WNCs with short packet transmissions," *IEEE Internet Things J.*, 2024.
- [56] F. Song, Q. Yang, M. Deng, H. Xing, Y. Liu, X. Yu, K. Li, and L. Xu, "AoI and energy tradeoff for aerial-ground collaborative MEC: A multi-objective learning approach," *IEEE Trans. Mob. Comput.*, vol. 23, no. 12, pp. 11 278–11 294, 2024.
- [57] M. Azizi, F. Zeinali, M. R. Mili, and S. Shokrollahi, "Efficient AoI-aware resource management in VLC-V2X networks via multi-agent RL mechanism," *IEEE Trans. Veh. Technol.*, vol. 73, no. 9, pp. 14 009–14 014, 2024.
- [58] S. Sun, W. Wu, C. Fu, X. Qiu, J. Luo, and J. Wang, "AoI optimization in multi-source update network systems under stochastic energy harvesting model," *IEEE J. Sel. Areas Commun.*, vol. 42, no. 11, pp. 3172–3187, 2024.
- [59] Q. Chen, S. Guo, W. Xu, J. Li, T. Shi, H. Gao, and Z. Cai, "Average AoI minimization with directional charging for wireless-powered network edge," *IEEE Transactions on Mobile Computing*, 2025.
- [60] X. Liu, H. Liu, K. Zheng, J. Liu, T. Taleb, and N. Shiratori, "AoI-minimal clustering, transmission and trajectory co-design for UAV-assisted WPCNs," *IEEE Trans. Veh. Technol.*, vol. 74, no. 1, pp. 1035–1051, 2025.

- [61] M. Yi, X. Wang, J. Liu, Y. Zhang, and R. Hou, "Multitask transfer deep reinforcement learning for timely data collection in rechargeable-UAV-aided IoT networks," *IEEE Internet Things J.*, vol. 10, no. 23, pp. 20 545–20 559, 2023.
- [62] X. Zhang, Z. Chang, T. Hämäläinen, and G. Min, "AoI-energy tradeoff for data collection in UAV-assisted wireless networks," *IEEE Trans. Commun.*, vol. 72, no. 3, pp. 1849–1861, 2024.
- [63] H. Zhao, G. Lu, Y. Liu, Z. Chang, L. Wang, and T. Hämäläinen, "Safe DQN-based AoI-minimal task offloading for UAV-aided edge computing system," *IEEE Internet Things J.*, vol. 11, no. 19, pp. 32 012–32 024, 2024.
- [64] J. Huang, T. Yu, F. Yang, S. Zhang, W. Jiang, and D. Niyato, "AoI-aware resource allocation with interference avoidance for ultra-dense industrial Internet of Things networks," *IEEE Internet Things J.*, 2024.
- [65] Y. Liu, X. Wang, G. Zheng, X. Wan, and Z. Ning, "An AoI-aware data transmission algorithm in blockchain-based intelligent healthcare systems," *IEEE Trans. Consumer Electron.*, vol. 70, no. 1, pp. 1180–1190, 2024.
- [66] C. Lin and W. Liao, "AoI-aware interference mitigation for task-oriented multicasting in multi-cell NOMA networks," *IEEE Trans. Wirel. Commun.*, vol. 23, no. 9, pp. 11 341–11 356, 2024.
- [67] Z. Huang, W. Wu, C. Fu, V. Chau, X. Liu, J. Wang, and J. Luo, "AoI-guaranteed bandit: Information gathering over unreliable channels," *IEEE Trans. Mob. Comput.*, vol. 23, no. 10, pp. 9469–9486, 2024.
- [68] Z. Li, F. Hu, Q. Li, Z. Ling, Z. Chang, and T. Hämäläinen, "AoI-aware waveform design for cooperative joint radar-communications systems with online prediction of radar target property," *IEEE Trans. Commun.*, vol. 72, no. 10, pp. 6029–6043, 2024.
- [69] K. Qi, Q. Wu, P. Fan, N. Cheng, W. Chen, J. Wang, and K. B. Letaief, "Deep-reinforcement-learning-based AoI-aware resource allocation for RIS-aided IoV networks," *IEEE Trans. Veh. Technol.*, vol. 74, no. 1, pp. 1365–1378, 2025.
- [70] X. Zhang, K. Xiong, W. Chen, P. Fan, B. Ai, and K. B. Letaief, "Minimizing AoI in high-speed railway mobile networks: DQN-based methods," *IEEE Trans. Intell. Transp. Syst.*, vol. 25, no. 12, pp. 20 137–20 150, 2024.
- [71] T. Shi, Q. Xu, J. Wang, C. Xu, K. Wu, K. Lu, and C. Qiao, "Enhancing the safety of autonomous driving systems via AoI-optimized task scheduling," *IEEE Trans. Veh. Technol.*, vol. 74, no. 3, pp. 3804–3819, 2025.
- [72] H. Xie, S.-W. Jeon, and H. Jin, "Distributed real-time control for minimizing AoI in random access networks," *IEEE Internet Things J.*, 2024.
- [73] Y. Qin, M. A. Kishk, and M. Alouini, "Velocity-aware statistical analysis of peak AoI for ground and aerial users," *IEEE Trans. Veh. Technol.*, 2025.
- [74] M. Akbari, A. Syed, W. S. Kennedy, and M. Erol-Kantarci, "AoI-aware energy-efficient SFC in UAV-aided smart agriculture using asynchronous federated learning," *IEEE Open J. Commun. Soc.*, vol. 5, pp. 1222–1242, 2024.
- [75] C. Liu, Y. Guo, N. Li, and X. Song, "AoI-minimal task assignment and trajectory optimization in multi-UAV-assisted IoT networks," *IEEE Internet Things J.*, vol. 9, no. 21, pp. 21 777–21 791, 2022.
- [76] J. Dai, Y. Gao, C. Cai, W. Xiong, and M. Liu, "UAV-enabled inspection system with no-fly zones: DRL-based joint mobile nest scheduling and UAV trajectory design," *IEEE Access*, vol. 13, pp. 10 844–10 856, 2025.
- [77] K. Meng, C. Chen, T. Wu, B. Xin, M. Liang, and F. Deng, "Evolutionary state estimation-based multi-strategy jellyfish search algorithm for multi-UAV cooperative path planning," *IEEE Trans. Intell. Veh.*, 2024.
- [78] C. Peng, X. Huang, Y. Wu, and J. Kang, "Constrained multi-objective optimization for UAV-enabled mobile edge computing: Offloading optimization and path planning," *IEEE Wirel. Commun. Lett.*, vol. 11, no. 4, pp. 861–865, 2022.
- [79] G. Sun, Y. Wang, Z. Sun, Q. Wu, J. Kang, D. Niyato, and V. C. M. Leung, "Multi-objective optimization for multi-UAV-assisted mobile edge computing," *IEEE Trans. Mob. Comput.*, vol. 23, no. 12, pp. 14 803–14 820, 2024.
- [80] Y. Zhou, A. A. Khuwaja, X. Li, N. Zhao, and Y. Chen, "Optimizing multi-UAV multi-user system through integrated sensing and communication for Age of Information (AoI) analysis," *IEEE Open J. Commun. Soc.*, vol. 5, pp. 6918–6931, 2024.
- [81] X. Gao, X. Zhu, and L. Zhai, "AoI-sensitive data collection in multi-UAV-assisted wireless sensor networks," *IEEE Trans. Wirel. Commun.*, vol. 22, no. 8, pp. 5185–5197, 2023.
- [82] S. A. Ullah, M. A. Sohail, H. Jung, M. O. B. Saeed, and S. A. Hassan, "Sum rate maximization in IoT networks with diversity-enhanced energy harvesting: A DRL-guided approach," *IEEE Internet Things J.*, vol. 11, no. 18, pp. 30 309–30 322, 2024.
- [83] J. Pan, Y. Li, R. Chai, S. Xia, and L. Zuo, "Multi-objective trajectory planning for UAV-assisted IoT networks based on DRL approach," *IEEE Internet Things J.*, 2025.
- [84] R. Zhang, H. Du, Y. Liu, D. Niyato, J. Kang, Z. Xiong, A. Jamalipour, and D. I. Kim, "Generative AI agents with large language model for satellite networks via a mixture of experts transmission," *IEEE J. Sel. Areas Commun.*, vol. 42, no. 12, pp. 3581–3596, 2024.
- [85] R. Zhang, K. Xiong, Y. Lu, P. Fan, D. W. K. Ng, and K. B. Letaief, "Energy efficiency maximization in RIS-assisted SWIPT networks with RSMA: A PPO-based approach," *IEEE J. Sel. Areas Commun.*, vol. 41, no. 5, pp. 1413–1430, 2023.
- [86] F. H. Panahi, F. H. Panahi, and T. Ohtsuki, "A reinforcement learning-based fire warning and suppression system using unmanned aerial vehicles," *IEEE Trans. Instrum. Meas.*, vol. 72, pp. 1–16, 2023.
- [87] O. S. Oubbati, A. Lakas, and M. Guizani, "Multiagent deep reinforcement learning for wireless-powered UAV networks," *IEEE Internet Things J.*, vol. 9, no. 17, pp. 16 044–16 059, 2022.
- [88] X. He, S. Pang, H. Gui, K. Zhang, N. Wang, and X. Zhai, "Multi-agent DRL-based large-scale heterogeneous task offloading for dynamic IoT systems," *IEEE Trans. Netw. Sci. Eng.*, vol. 12, no. 2, pp. 982–996, 2025.
- [89] W. Wang, X. Xu, M. Bilal, M. Khan, and Y. Xing, "UAV-assisted content caching for human-centric consumer applications in IoV," *IEEE Trans. Consumer Electron.*, vol. 70, no. 1, pp. 927–938, 2024.
- [90] P. Qin, Y. Fu, J. Zhang, S. Geng, J. Liu, and X. Zhao, "DRL-based resource allocation and trajectory planning for NOMA-enabled multi-UAV collaborative caching 6G network," *IEEE Trans. Veh. Technol.*, vol. 73, no. 6, pp. 8750–8764, 2024.
- [91] B. Yin, X. Fang, and X. Wang, "Joint optimization of trajectory control, resource allocation, and user association based on DRL for multi-fixed-wing UAV networks," *IEEE Trans. Wirel. Commun.*, vol. 23, no. 10, pp. 13 330–13 343, 2024.
- [92] F. Li, C. Gu, D. Liu, Y. Wu, and H. Wang, "DRL-based joint task scheduling and trajectory planning method for UAV-assisted MEC scenarios," *IEEE Access*, vol. 12, pp. 156 224–156 234, 2024.
- [93] G. B. Tarekegn, R.-T. Juang, B. A. Tesfaw, H.-P. Lin, H.-C. Hsu, R. B. Tarekegn, and L.-C. Tai, "A centralized multi-agent DRL-based trajectory control strategy for unmanned aerial vehicle-enabled wireless communications," *IEEE Open J. Veh. Technol.*, 2024.
- [94] Y. Wei, Y. Lu, P. Zhao, S. Leng, and K. Yang, "Minimizing age of information in UAV-assisted data collection with limited charging facilities," *IEEE Wirel. Commun. Lett.*, vol. 13, no. 5, pp. 1463–1467, 2024.
- [95] K. Messaoudi, A. Baz, O. S. Oubbati, A. Rachedi, T. Bendouma, and M. Atiquzzaman, "UGV charging stations for UAV-assisted AoI-aware data collection," *IEEE Trans. Cogn. Commun. Netw.*, vol. 10, no. 6, pp. 2325–2343, 2024.
- [96] M. L. Betalo, S. Leng, H. N. Abishu, A. M. Seid, M. Fakirah, A. Erbad, and M. Guizani, "Multi-agent DRL-based energy harvesting for freshness of data in UAV-assisted wireless sensor networks," *IEEE Trans. Netw. Serv. Manag.*, vol. 21, no. 6, pp. 6527–6541, 2024.
- [97] J. Li, H. Zhao, H. Wang, F. Gu, J. Wei, H. Yin, and B. Ren, "Joint optimization on trajectory, altitude, velocity, and link scheduling for minimum mission time in UAV-aided data collection," *IEEE Internet Things J.*, vol. 7, no. 2, pp. 1464–1475, 2020.
- [98] Y. Wu and K. H. Low, "Route coordination of UAV fleet to track a ground moving target in search and lock (SAL) task over urban airspace," *IEEE Internet Things J.*, vol. 9, no. 20, pp. 20 604–20 619, 2022.
- [99] H. Wu, F. Lyu, C. Zhou, J. Chen, L. Wang, and X. Shen, "Optimal UAV caching and trajectory in aerial-assisted vehicular networks: A learning-based approach," *IEEE J. Sel. Areas Commun.*, vol. 38, no. 12, pp. 2783–2797, 2020.
- [100] J. Li, G. Sun, X. Sun, F. Mei, J. Wang, X. Hou, D. Tian, and V. C. M. Leung, "Securing the sky: Integrated satellite-UAV physical layer security for low-altitude wireless networks," *arXiv preprint arXiv:2506.23493*, 2025.
- [101] S. Zhang, W. Liu, and N. Ansari, "Joint wireless charging and data collection for UAV-enabled Internet of Things network," *IEEE Internet Things J.*, vol. 9, no. 23, pp. 23 852–23 859, 2022.
- [102] Z. Wang, R. Liu, Q. Liu, J. S. Thompson, and M. Kadoch, "Energy-efficient data collection and device positioning in UAV-assisted IoT," *IEEE Internet Things J.*, vol. 7, no. 2, pp. 1122–1139, 2020.
- [103] H. Lu, F. Lyu, J. Ren, H. Wu, C. Zhou, Z. Liu, Y. Zhang, and X. Shen, "CODE⁺: Fast and accurate inference for compact distributed IoT

- data collection," *IEEE Trans. Parallel Distributed Syst.*, vol. 35, no. 11, pp. 2006–2022, 2024.
- [104] W. Yuan, Y. Cui, J. Wang, F. Liu, G. Sun, T. Xiang, J. Xu, S. Jin, D. Niyato, S. Coleri, S. Sun, S. Mao, A. Jamalipour, D. I. Kim, M.-S. Alouini, and X. Shen, "From ground to sky: Architectures, applications, and challenges shaping low-altitude wireless networks," *arXiv preprint arXiv:2506.12308*, 2025.
- [105] Z. Wei, M. Zhu, N. Zhang, L. Wang, Y. Zou, Z. Meng, H. Wu, and Z. Feng, "UAV-assisted data collection for Internet of Things: A survey," *IEEE Internet Things J.*, vol. 9, no. 17, pp. 15 460–15 483, 2022.
- [106] X. Xu, H. Zhao, H. Yao, and S. Wang, "A blockchain-enabled energy-efficient data collection system for UAV-assisted IoT," *IEEE Internet Things J.*, vol. 8, no. 4, pp. 2431–2443, 2021.
- [107] H. Lu, F. Lyu, H. Wu, J. Zhang, J. Ren, Y. Zhang, and X. Shen, "FL-AMM: federated learning augmented map matching with heterogeneous cellular moving trajectories," *IEEE J. Sel. Areas Commun.*, vol. 41, no. 12, pp. 3878–3892, 2023.
- [108] S. Duan, F. Lyu, J. Zhang, H. Lu, P. Yang, H. Wu, Y. Zhang, and X. Shen, "MoCo: Urban user mobile contact detection based on cellular signaling trace," *IEEE Trans. Mob. Comput.*, 2025.
- [109] Q. Zhang, W. Fang, Q. Liu, J. Wu, P. Xia, and L. Yang, "Distributed laser charging: A wireless power transfer approach," *IEEE Internet Things J.*, vol. 5, no. 5, pp. 3853–3864, 2018.
- [110] W. Jaafar and H. Yanikomeroglu, "Dynamics of laser-charged UAVs: A battery perspective," *IEEE Internet Things J.*, vol. 8, no. 13, pp. 10 573–10 582, 2021.
- [111] M. Zhao, Q. Shi, and M. Zhao, "Efficiency maximization for UAV-enabled mobile relaying systems with laser charging," *IEEE Trans. Wirel. Commun.*, vol. 19, no. 5, pp. 3257–3272, 2020.
- [112] T. T. Nguyen, N. D. Nguyen, and S. Nahavandi, "Deep reinforcement learning for multiagent systems: A review of challenges, solutions, and applications," *IEEE Trans. Cybern.*, vol. 50, no. 9, pp. 3826–3839, 2020.
- [113] C. Guo, X. Wang, Y. Zheng, and F. Zhang, "Real-time optimal energy management of microgrid with uncertainties based on deep reinforcement learning," *Energy*, vol. 238, p. 121873, 2022.
- [114] S. Duan, F. Lyu, H. Wu, W. Chen, H. Lu, Z. Dong, and X. Shen, "MOTO: mobility-aware online task offloading with adaptive load balancing in small-cell MEC," *IEEE Trans. Mob. Comput.*, vol. 23, no. 1, pp. 645–659, 2024.
- [115] Y. Deng, F. Lyu, J. Ren, H. Wu, Y. Zhou, Y. Zhang, and X. Shen, "AUCTION: automated and quality-aware client selection framework for efficient federated learning," *IEEE Trans. Parallel Distributed Syst.*, vol. 33, no. 8, pp. 1996–2009, 2022.
- [116] J. Nievergelt, "Exhaustive search, combinatorial optimization and enumeration: Exploring the potential of raw computing power," in *Proc. Springer SOFSEM*, vol. 1963, 2000, pp. 18–35.
- [117] S. P. Boyd and L. Vandenberghe, *Convex optimization*. Cambridge university press, 2004.
- [118] C. A. Bliss, M. R. Frank, C. M. Danforth, and P. S. Dodds, "An evolutionary algorithm approach to link prediction in dynamic social networks," *J. Comput. Sci.*, vol. 5, no. 5, pp. 750–764, 2014.
- [119] I. K. Nikolos, K. P. Valavanis, N. Tsourveloudis, and A. N. Kostaras, "Evolutionary algorithm based offline/online path planner for UAV navigation," *IEEE Trans. Syst. Man Cybern. Part B*, vol. 33, no. 6, pp. 898–912, 2003.
- [120] W. Zhang, Q. Wang, X. Liu, Y. Liu, and Y. Chen, "Three-dimension trajectory design for multi-UAV wireless network with deep reinforcement learning," *IEEE Trans. Veh. Technol.*, vol. 70, no. 1, pp. 600–612, 2021.
- [121] H. Mei, K. Yang, Q. Liu, and K. Wang, "3D-trajectory and phase-shift design for RIS-assisted UAV systems using deep reinforcement learning," *IEEE Trans. Veh. Technol.*, vol. 71, no. 3, pp. 3020–3029, 2022.
- [122] H. Bayerlein, M. Theile, M. Caccamo, and D. Gesbert, "Multi-UAV path planning for wireless data harvesting with deep reinforcement learning," *IEEE Open J. Commun. Soc.*, vol. 2, pp. 1171–1187, 2021.
- [123] B. Li and Y. Wu, "Path planning for UAV ground target tracking via deep reinforcement learning," *IEEE Access*, vol. 8, pp. 29 064–29 074, 2020.
- [124] C. Yu, A. Velu, E. Vinitsky, J. Gao, Y. Wang, A. M. Bayen, and Y. Wu, "The surprising effectiveness of PPO in cooperative multi-agent games," in *Proc. NeurIPS*, 2022, pp. 24 611–24 624.
- [125] J. Ackermann, V. Gabler, T. Osa, and M. Sugiyama, "Reducing overestimation bias in multi-agent domains using double centralized critics," *arXiv preprint arXiv:1910.01465*, 2019.
- [126] R. Lowe, Y. Wu, A. Tamar, J. Harb, P. Abbeel, and I. Mordatch, "Multi-agent actor-critic for mixed cooperative-competitive environments," in *Proc. NeurIPS*, 2017, pp. 6379–6390.
- [127] S. Hahm, J. Kim, A. Jeong, H. Yi, S. Chang, S. N. Kishore, A. Chauhan, and S. P. Cherian, "Reliable real-time operating system for IoT devices," *IEEE Internet Things J.*, vol. 8, no. 5, pp. 3705–3716, 2021.
- [128] M. Bukhsh, S. Abdullah, A. Rahman, M. N. Asghar, H. Arshad, and A. Alabdulatif, "An energy-aware, highly available, and fault-tolerant method for reliable IoT systems," *IEEE Access*, vol. 9, pp. 145 363–145 381, 2021.
- [129] Z. Xia, J. Du, J. Wang, C. Jiang, Y. Ren, G. Li, and Z. Han, "Multi-agent reinforcement learning aided intelligent UAV swarm for target tracking," *IEEE Trans. Veh. Technol.*, vol. 71, no. 1, pp. 931–945, 2022.
- [130] P. B. Sujit and J. B. de Sousa, "Multi-UAV task allocation with communication faults," in *Proc. IEEE ACC*, 2012, pp. 3724–3729.
- [131] A. Tahir, M. H. Haghbayan, J. M. Böling, and J. Plosila, "Energy-efficient post-failure reconfiguration of swarms of unmanned aerial vehicles," *IEEE Access*, vol. 11, pp. 24 768–24 779, 2023.
- [132] L. Xing and B. W. Johnson, "Reliability theory and practice for unmanned aerial vehicles," *IEEE Internet Things J.*, vol. 10, no. 4, pp. 3548–3566, 2023.
- [133] K. R. Weiss, T. M. Khoshgoftaar, and D. Wang, "A survey of transfer learning," *J. Big Data*, vol. 3, p. 9, 2016.
- [134] L. Zhao, Y. Han, A. Hawbani, S. Wan, Z. Guo, and M. Guizani, "MEDIA: an incremental DNN based computation offloading for collaborative cloud-edge computing," *IEEE Trans. Netw. Sci. Eng.*, vol. 11, no. 2, pp. 1986–1998, 2024.