

Graphical Abstract

DRL-based Task Offloading for Networked UAVs with Random Mobility and Collision Avoidance*

Xixin Zhang, Dongge Jia, Junfei Xie

Highlights

DRL-based Task Offloading for Networked UAVs with Random Mobility and Collision Avoidance

Xixin Zhang, Dongge Jia, Junfei Xie

- Research highlight 1
- Research highlight 2

DRL-based Task Offloading for Networked UAVs with Random Mobility and Collision Avoidance

Xixin Zhang^{a,b,1}, Dongge Jia^{b,2}, Junfei Xie^{b,3,*}

^a*Department of Electrical and Computer Engineering, University of California San Diego, 9500 Gilman Dr, La Jolla, 92092, California, USA*

^b*Department of Electrical and Computer Engineering, San Diego State University, 5500 Campanile Drive, San Diego, 92182, California, USA*

Abstract

Unmanned Aerial Vehicles (UAVs) have gained widespread use across various fields due to their flexibility and multifunctionality. However, their limited onboard computing capacity is often criticized for hindering their ability to execute complex tasks in real-time. To address this challenge, Networked Airborne Computing (NAC) has emerged, which leverages the collective computing power of multiple UAVs to enable efficient handling of large-scale data processing, real-time analytics, and complex mission coordination. Despite its potential, research in this area is still in its infancy. In this paper, we consider a typical NAC scenario where multiple UAVs with collision avoidance capabilities share resources while moving randomly within an area. Without prior knowledge of the system models, we aim to optimize task allocation among UAVs with uncertain mobility. To achieve this, we propose a Deep Reinforcement Learning algorithm based on the Twin Delayed Deep Deterministic Policy Gradient (TD3). Simulation results demonstrate that our approach significantly speeds up task execution compared to existing meth-

*This document is the results of the research project funded by the National Science Foundation.

^{*}Corresponding author

Email addresses: xiz166@ucsd.edu (Xixin Zhang), xxxx@email.edu (Dongge Jia), jxie4@sdsu.edu (Junfei Xie)

¹This is the first author footnote.

²Another author footnote, this is a very long footnote and it should be a long footnote. But this footnote is not yet sufficiently long enough to make two lines of footnote text.

³Yet another author footnote.

ods.

Keywords: Computation Offloading, Deep Reinforcement Learning, Unmanned Aerial Vehicle, Edge Computing

1. Introduction

In recent years, unmanned aerial vehicles (UAVs), or drones, have seen rapid advancements and growing popularity in areas such as precision agriculture, disaster response, aerial photography, and environmental monitoring [1, 2]. As UAV applications become increasingly complex, the use of multiple cooperative UAVs has become more common. Nevertheless, their limited onboard computational resources often become a bottleneck. One solution that naturally follows is to offload computationally intensive tasks to external resources.

Extensive research has focused on efficiently utilizing resources on edge servers or remote clouds to support multi-UAV applications. For instance, Liu *et al.* [3] proposed to utilize a UAV-Edge-Cloud computing model and formulate a joint optimization of workflow assignment and multi-hop routing scheduling for UAV swarms to minimize computation cost and latency. Bai *et al.* [4] investigated delay-aware cooperative task offloading for multi-UAV enabled edge-cloud computing, proposing an algorithm to balance task distribution and minimize completion delay. In these studies, UAVs in the swarm are typically viewed as relays that bring edge servers or remote clouds closer, rather than computing nodes. They get sufficient computing resources at the cost of a high data transmission delay, which may not be acceptable for time-sensitive UAV applications not to mention real-time tasks. Moreover, mobile edge servers require a reliable local network infrastructure, which is difficult to deploy and scale, especially in underdeveloped or post-disaster areas [5].

With technological advancements, the emergence of small, lightweight yet powerful micro-computers has significantly accelerated the onboard computing capacity of UAVs. This has spurred researchers to explore UAVs' potential in acting as edge servers. In [6], Hu *et al.* leveraged the computing resources of a moving UAV to serve ground users, aimed to minimize the total maximum delays among users by jointly optimizing offloading ratios, user scheduling, and UAV trajectory in a UAV-aided mobile edge computing system. Miao *et al.* [7] proposed a multi-UAV-assisted mobile edge com-

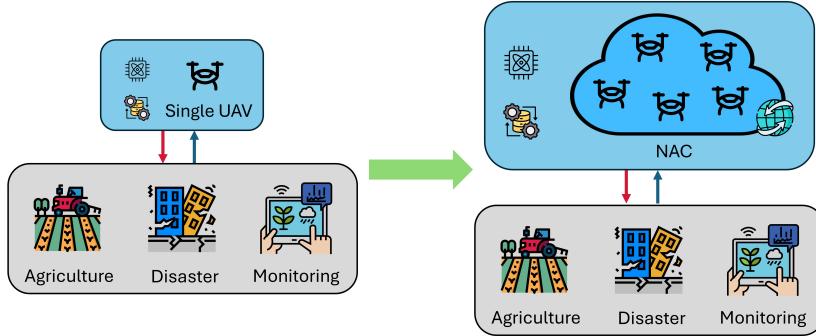


Figure 1: Networked Airborne Computing (NAC)

33 putting (MEC) offloading algorithm that maximizes the access quantity and
 34 minimizes the task completion latency by cluster path planning based on
 35 user mobility and communication coverage. Although UAVs have proven
 36 promising in providing on-demand computing resources, these studies treat
 37 them as separate servers.

38 To harness the full computational potential of multi-UAV systems, a new
 39 paradigm called Networked Airborne Computing (NAC) is proposed, where
 40 multiple aerial vehicles share resources among each other[8] as in Fig.??.
 41 The fast deployment, infrastructure-free, and low-cost characteristics make
 42 the UAV-based NAC a promising technique. Nevertheless, research in NAC
 43 is still in its early stages. In our previous studies, we have developed a ROS-
 44 based simulator and a hardware testbed that consists of multiple UAVs to
 45 facilitate NAC research [9]. In [10], we introduced a coded distributed com-
 46 puting scheme based on deep reinforcement learning (DRL) for optimally
 47 partitioning and allocating tasks to multiple networked UAVs. This scheme
 48 addresses two typical NAC scenarios. The first scenario involves uncontrol-
 49 lable UAV mobility, which can happen when they are operated by different
 50 owners. In the second scenario, UAVs are controlled to assist in task com-
 51 putation. Simulation results demonstrate the effectiveness of the proposed
 52 scheme. However, in the first scenario, we assumed UAVs maintain a consis-
 53 tent movement pattern throughout the execution of a particular task and did
 54 not account for motion interference between UAVs due to collision avoidance.
 55 Moreover, the simple matrix multiplication tasks were considered.

56 Orthogonal Frequency Division Multiple Access (OFDMA) is a sophis-
 57 ticated wireless communication technology that divides the available band-

58 width into multiple orthogonal subcarriers, dynamically allocating them to
59 users based on their channel conditions and service requirements [11]. This
60 dynamic resource allocation is particularly advantageous in networked Un-
61 manned Aerial Vehicle (UAV) systems, where the mobility and dynamic
62 topology of UAVs demand robust and flexible communication solutions. By
63 leveraging OFDMA, networked UAV systems can achieve high spectral ef-
64 ficiency, reduced latency, and reliable performance in diverse environments,
65 supporting applications such as real-time surveillance, package delivery, and
66 disaster response [12]. Energy efficiency is another critical concern in UAV
67 systems due to the limited onboard battery capacity, which constrains mis-
68 sion duration and operational effectiveness. Task offloading, while reduc-
69 ing onboard computational load, introduces additional energy costs for data
70 transmission and reception. This makes it essential to adopt optimization
71 strategies that minimize total energy consumption—an especially pressing
72 challenge for energy-limited systems such as the NAC system. In addressing
73 these challenges, [13] proposes an optimization framework that simultane-
74 ously allocates subcarriers and adjusts power levels to balance spectral ef-
75 ficiency and energy consumption. This approach is particularly effective in
76 dynamic multi-UAV communication networks, where varying channel con-
77 ditions and interference necessitate adaptive and efficient resource manage-
78 ment.

79 In this paper, we investigate a more common yet challenging NAC sce-
80 nario where all UAVs, including both offloaders and offloadees, move ran-
81 domly during task execution while actively avoiding collisions. None of the
82 UAVs have prior knowledge of the environment or system models, and their
83 movement patterns or future trajectories are not shared among each other.
84 Additionally, we generalize computation tasks as any functions or operations
85 that can be partitioned into arbitrary subtasks for parallel computation. To
86 model UAV movement, we extend the traditional Random Direction model
87 [14], originally designed for individual entities, to capture collision avoidance
88 interactions among multiple UAVs. Furthermore, we formulate a nonlin-
89 ear optimization to optimize task allocation and develop a DRL algorithm
90 based on the Twin Delayed Deep Deterministic Policy Gradient (TD3)[15]
91 to solve it. We evaluate the performance of the proposed method through
92 extensive comparative simulation studies, which demonstrate its promising
93 performance.

94 In the rest of this paper, Sec. 3 details the system models and formulates
95 the optimization problem. Sec. 5 describes the proposed DRL algorithm.

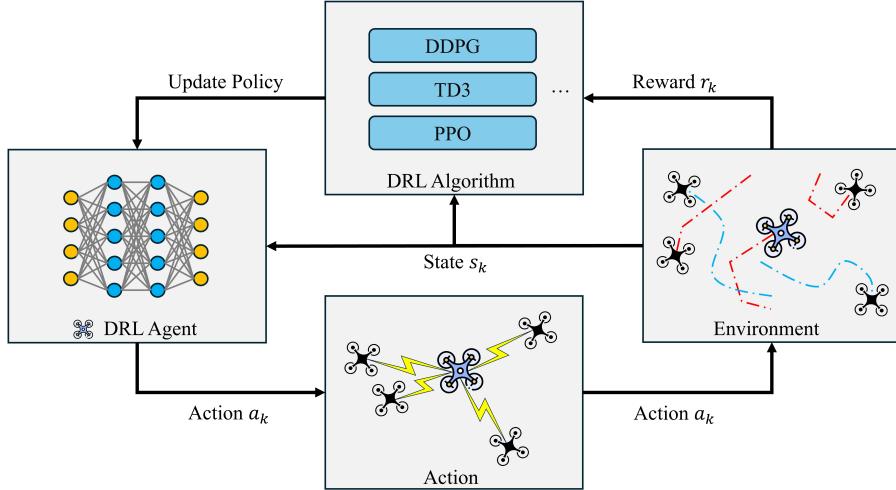


Figure 2: Caption

⁹⁶ In Sec. 6, simulation results are presented and discussed. We conclude in
⁹⁷ Sec. ??.

⁹⁸ 2. Related Work

⁹⁹ 3. System Models

¹⁰⁰ In this section, we will introduce the system models. Consider a group of
¹⁰¹ $N + 1$ heterogeneous UAVs with varying physical configurations, indexed as
¹⁰² $i \in \mathcal{N} = \{0, 1, 2, \dots, N\}$. Each UAV is equipped with computing and commu-
¹⁰³ nication modules, enabling resource sharing and onboard computation. Their
¹⁰⁴ characteristics of computing, communication, energy consumption, and mo-
¹⁰⁵ bility within the system can be comprehensively described and modeled as
¹⁰⁶ follows

¹⁰⁷ 3.1. Computing Model

¹⁰⁸ We describe the computing capability of each UAV i as idle CPU cycle
¹⁰⁹ frequency $f_i(t)$ in Hz, which is variable dependent on time and can fluctuate
¹¹⁰ due to some ongoing computing tasks. For a general computing task k ,
¹¹¹ its input data size is S_k (bits), and its required computation intensity is ξ_k

112 (cycles/bit)[16]. The total CPU cycles required to compute task k is hence
113 $\xi_k S_k$ and the time required for UAV i to execute this task is

$$T_{k,i}^{comp} = \frac{\xi_k \cdot S_k}{f_i(t)} \quad (1)$$

114 The corresponding energy consumption during computing can be given as

$$E_{k,i}^{comp} = \epsilon \cdot f_i(t)^3 \cdot T_{k,i}^{comp} \quad (2)$$

115 where ϵ represents an energy consumption parameter associated with the
116 effective switched capacitance, which is determined by the underlying CPU
117 architecture. To simplify the analysis, it is assumed that this capacitance
118 remains uniform across all devices [17].

119 3.2. Communication Model

120 At time t , let the distance between UAV i and UAV j be denoted as $d_{ij}(t)$.
121 The UAV-to-UAV links are typically Line of Sight (LoS), with propagation
122 speed approaching the speed of light. Hence, the transmission latency can be
123 approximated using the transmission time. Here, we model the transmission
124 rate (bits/s) based on the Simplified Path Loss Model [18] as follows

$$\nu_{ij}(t) = \begin{cases} B_{ij}(t) \log_2 \left(1 + \frac{G(d_r/d_{ij}(t))^{\theta} \psi_{ij}(t)}{N_0 B_{ij}(t)} \right), & \text{if } \nu_{ij} \geq \nu_{\min}, \\ 0, & \text{otherwise.} \end{cases} \quad (3)$$

125 where $B_{ij}(t)$ is the bandwidth (Hz) of the specific channel between UAV i
126 and j , d_r is constant reference distance (meter), G is the unitless constant
127 equal to the path gain of the distance d_r , θ is the path loss exponent, $\psi_{ij}(t)$
128 is the transmitted power (mW) and N_0 is the constant noise power spectral
129 density (dBm/Hz), ν_{ij} is the threshold to suppress the data rates that are
130 too low, which can be regarded as the constraint on QoS.

131 Unlike our previous work [19] which assumed the channel bandwidth and
132 transmitted power as constants, here we treat both B_{ij} and ψ_{ij} as variables
133 to be determined for the specific task k . Once the bandwidth B_{ij}^k and trans-
134 mission power ψ_{ij}^k allocated for task k are determined, the data rate ν_{ij}^k is also
135 determined follows Eq.3. Then the overall transmission time is as follows

$$T_{k,ij}^{trans} = \frac{S_k}{\nu_{ij}^k} \quad (4)$$

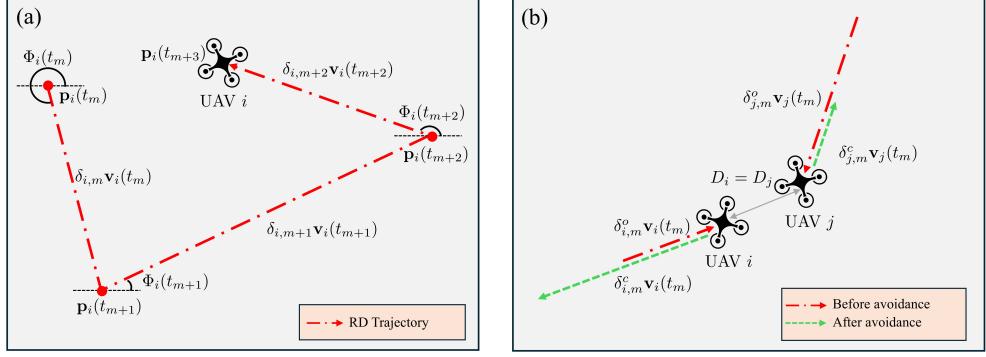


Figure 3: Caption

136 The reception power of UAV j is denoted as $\tilde{\psi}_j$ (mW) for completeness as
 137 in [20], then the energy (Joule) consumed for offloading task by the system
 138 is a combination of transmission energy and reception energy as

$$E_{k,ij}^{trans} = (\psi_{ij} + \tilde{\psi}_j) \cdot 10^{-3} \cdot T_{k,ij}^{trans} \quad (5)$$

139 *3.3. Mobility Model*

140 We assume the UAVs fly at the same altitude. Therefore, the position
 141 of each UAV i at time t can be depicted as $\mathbf{p}_i(t) = (x_i(t), y_i(t)) \in \mathbb{R}^2$ with
 142 constraints $0 \leq x_i(t) \leq W$, $0 \leq y_i(t) \leq W$, such that the position is bounded
 143 within an area of $W \times W(m^2)$. Given initial position $\mathbf{p}_i(0) = (x_i(0), y_i(0))$,
 144 we adapt the Random Direction (RD) model [14] and Smooth Turn (ST)
 145 model[21] to model its movement, which have been widely used for describing
 146 UAVs, particularly multirotor drones [22]. At time t , let the velocity of UAV
 147 i be $\mathbf{v}_i(t) = (v_{i,x}(t), v_{i,y}(t)) \in \mathbb{R}^2$ and the heading direction be $\Phi_i(t) \in [0, 2\pi]$,
 148 where $v_{i,x}(t)$, $v_{i,y}(t)$ are component of velocity in x and y direction.

149 *3.3.1. Random Direction*

150 If UAV i moves in the mode of Random Direction (RD), it randomly
 151 picks a velocity and moves along a straight line at this velocity for a duration
 152 randomly picked as well until another set of velocity and duration is selected
 153 and repeats the process (Fig 3a).

154 Suppose the start time of m -th duration is denoted as $t_m = \sum_{l=0}^{m-1} \delta_{i,l}$
 155 which is also the m -th instance when UAV i changes its velocity, where each
 156 $\delta_{i,l}$ is randomly sampled from exponential distribution with parameter λ_{rd}

157 and t_0 is set to be 0. At time t_m , UAV i select a speed magnitude uniformly
 158 between 0 and $v_{max} \in \mathbb{R}$, i.e. $0 < |\mathbf{v}_i(t_m)| < v_{max}$ and the heading direction
 159 $\Phi_i(t_m)$ in radian uniformly across 2π , leading to the new velocity $\mathbf{v}_i(t_m)$ for
 160 the m -th duration such that

$$\begin{aligned} v_{i,x}(t_m) &= |\mathbf{v}_i(t_m)| \cos(\Phi_i(t_m)) \\ v_{i,y}(t_m) &= |\mathbf{v}_i(t_m)| \sin(\Phi_i(t_m)) \end{aligned}$$

161 As velocity remains unchanged within each duration, the UAV i 's position
 162 at time t , where $t_m \leq t < t_{m+1}$, can be represented as follows

$$\mathbf{p}_i(t) = \mathbf{p}_i(0) + \sum_{l=0}^{m-1} \delta_{i,l} \mathbf{v}_i(t_l) + (t - t_m) \mathbf{v}_i(t_m) \quad (6)$$

163 The traditional RD model [14] is originally designed to describe the mo-
 164 bility of independent single entities, which ignores the spatiotemporal cor-
 165 relations of trajectory across entities. However, in multi-UAV systems, the
 166 mobility of UAVs can change to avoid collisions. This necessitates the in-
 167 corporation of collision avoidance mechanisms into the RD model. Here, we
 168 assume that each UAV i will turn around and move in the opposite direction
 169 without changing speed until the current duration is completed when its dis-
 170 tance to any other UAV j falls below a threshold D_i . Note that if both UAVs
 171 have the same threshold $D_i = D_j$, UAV j will also reserve its direction (Fig.
 172 3b). By treating the boundaries of the area as obstacles, we ensure that all
 173 UAVs move within the designated area. The mobility of each UAV i with
 174 collision avoidance can then be described as

$$\begin{aligned} \mathbf{p}_i(t) = \mathbf{p}_i(0) &+ \sum_{l=0}^{m-1} (\delta_{i,l}^o - \delta_{i,l}^c) \mathbf{v}_i(t_l) \\ &+ (\tilde{\delta}_{i,m}^o - \tilde{\delta}_{i,m}^c) \mathbf{v}_i(t_m) \end{aligned} \quad (7)$$

175 where $0 \leq \delta_{i,l}^o \leq \delta_{i,l}$ is the time spent moving at velocity $\mathbf{v}_i(t_l)$ in the l -th
 176 instance before triggering collision avoidance and $\delta_{i,l}^c = \delta_{i,l} - \delta_{i,l}^o$. Likewise,
 177 $0 \leq \tilde{\delta}_{i,m}^o \leq t - t_m$ is the time spent moving at velocity $\mathbf{v}_i(t_m)$ before collision
 178 avoidance in the m -th instance, and $\tilde{\delta}_{i,m}^c = (t - t_m) - \tilde{\delta}_{i,m}^o$.

179 3.3.2. Smooth Turn

180 In the smooth turn (ST) mobility model, UAV i randomly picks a turning
 181 center located at a point along the line perpendicular to its current head-
 182 ing direction and moves at a constant forward speed following the circular

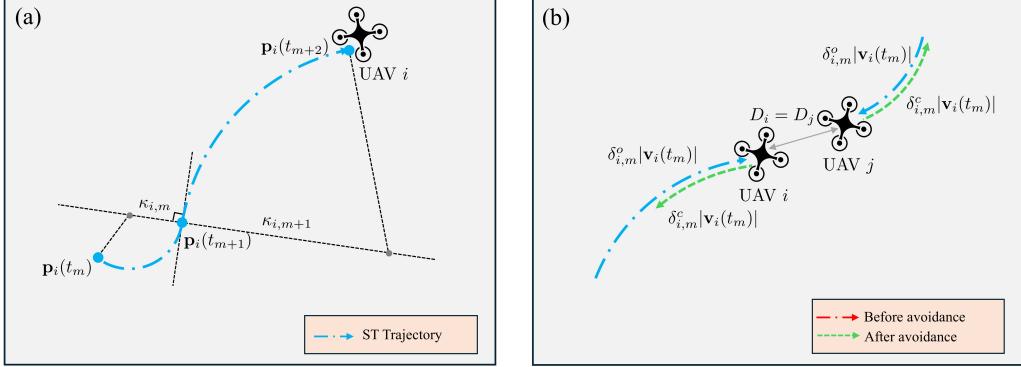


Figure 4: Caption

183 trajectory around the turning center for a duration randomly selected until
 184 another turning center picked and repeat the process.

185 Given the velocity magnitude $|\mathbf{v}_i(t)|$ constant for all t , UAV i changes its
 186 turning center at the start time $t_m = \sum_{l=0}^{m-1} \delta_{i,l}$ of the m -th duration, where
 187 each $\delta_{i,l}$ is randomly sampled from exponential distribution with parameter
 188 λ_{st} and t_0 is set to be 0. At time t_m , UAV i samples a variable $\kappa_{i,m} \in \mathbb{R}$
 189 from a Gaussian distribution as $\frac{1}{\kappa_{i,m}} \sim \mathcal{N}(0, \sigma^2)$. We define $|\kappa_{i,m}|$ as the
 190 radius of the circular trajectory, and the center of the circle is thereafter
 191 uniquely determined along the line perpendicular to the heading angle $\Phi_i(t_m)$,
 192 assuming counterclockwise rotation when $\kappa_{i,m} < 0$ and clockwise rotation
 193 when $\kappa_{i,m} < 0$. As $\kappa_{i,m} \in \mathbb{R}$ remains unchanged for the m -th duration, the
 194 heading angle $\Phi_i(t)$ of UAV i at time t , where $t_m \leq t < t_{m+1}$, gives as

$$\Phi_i(t) = \Phi_i(t_m) - \frac{|\mathbf{v}_i(t)|}{\kappa_{i,m}}(t - t_m) \quad (8)$$

195 and the velocity follows as

$$\begin{aligned} v_{i,x}(t) &= |\mathbf{v}_i(t)| \cos(\Phi_i(t)) \\ v_{i,y}(t) &= |\mathbf{v}_i(t)| \sin(\Phi_i(t)) \end{aligned}$$

196 The UAV i 's position at time t can be computed with $\mathbf{v}_i(t) = [v_{i,x}, v_{i,y}]$ as
 197 following

$$\mathbf{p}_i(t) = \mathbf{p}_i(0) + \sum_{l=0}^{m-1} \int_{t_l}^{t_{l+1}} \mathbf{v}_i(\tau) d\tau + \int_{t_m}^t \mathbf{v}_i(\tau) d\tau \quad (9)$$

198 The same issue appears in the original ST mobility model as the RD
 199 mobility model that it fails to model the interactions between nodes since it
 200 did not include collision avoidance mechanisms either. Here, we propose a
 201 collision avoidance strategy similar to the one introduced for the RD model
 202 in the previous section. In this approach, UAVs retrace their prior circular
 203 trajectory in the reverse direction, restoring the safe inter-UAV spacing and
 204 spatial relationship (Fig. 4b). The mobility in Eq. 10 of UAV i will become

$$\begin{aligned}
 \mathbf{p}_i(t) = & \mathbf{p}_i(0) + \sum_{l=0}^{m-1} \int_{t_l}^{t_l + \delta_{i,l}^o - \delta_{i,l}^c} \mathbf{v}_i(\tau) d\tau \\
 & + \int_{t_m}^{t_m + \tilde{\delta}_{i,m}^o - \tilde{\delta}_{i,m}^c} \mathbf{v}_i(\tau) d\tau
 \end{aligned} \tag{10}$$

205 where $\delta_{i,l}^o$ represents the duration that UAC i moves along the designated
 206 arc trajectory in its original direction the (determined by the sign of $\kappa_{i,l}$), and
 207 $\delta_{i,l}^c$ denotes the duration that UAV i moves in the opposite direction along
 208 the circular trajectory due to the collision avoidance maneuver within the
 209 l -th duration, such that $\delta_{i,l}^o + \delta_{i,l}^c = \delta_{i,m}$. Let $\tilde{\delta}_{i,m}$ be the duration in m -th
 210 instance until t , the splits of duration for the original direction and reversed
 211 direction are $\tilde{\delta}_{i,m}^o$, $\tilde{\delta}_{i,m}^c$ respectively so that $\tilde{\delta}_{i,m} = \tilde{\delta}_{i,m}^o + \tilde{\delta}_{i,m}^c$

212 4. Problem Formulation

213 Without loss of generality, we let UAV $i = 0$ be the master (or offloader)
 214 and treat the remaining N UAVs as potential offloadees with idle computing
 215 resources. Suppose a sequence of computing tasks $\mathcal{K} = \{1, 2, \dots, K\}$ is gener-
 216 ated at the master, and each task k can be divided into $L_k \in \mathbb{Z}^+$ atomic tasks
 217 of size $\ell_k = \frac{S_k}{L_k}$, which can be computed in parallel. Consider a UAV net-
 218 work utilizing OFDMA, where the total available communication bandwidth
 219 is represented by B . This bandwidth is evenly divided into W orthogonal
 220 subcarriers, with each subcarrier having a bandwidth of $b = \frac{B}{W}$. The master
 221 UAV, tasked with managing communication and data transmission, operates
 222 within a total power budget denoted by $\Psi \in \mathbb{R}$.

223 4.1. Joint Optimization

224 The goal of the master UAV is to minimize the total task completion time
 225 while ensuring minimal energy consumption. To achieve this, it strategically

226 allocates the L_k atomic tasks across all available UAVs, including itself. Furthermore, the master UAV dynamically assigns network resources to balance computational and communication loads, thereby ensuring efficient utilization of the network's bandwidth and power resources.

230 Suppose the workload offloaded to UAV $j \in \mathcal{N}$ for task k is $c_j^k \ell_k$, where c_j^k
231 is a non-negative integer that satisfies $0 \leq c_j^k \leq L_k$. Therefore, $\sum_{j=0}^N c_j^k = L_k$.
232 Of note, when there is no workload offloaded to UAV j , then $c_j^k = 0$. With
233 respect to the network resources, let the number of subcarriers allocated to
234 the channel between master and UAV $j \in \mathcal{N} \setminus \{0\}$ be non-negative integer
235 $0 \leq w_j^k \leq W$, the channel bandwidth for will be $B_{0j}^k = w_j^k b$. On the other
236 hand, the corresponding transmission power master UAV allocate to B_{0j}^k will
237 be $\psi_{0j}^k \in \mathbb{R}$ such that $\sum_{j=1}^N \psi_{0j}^k = \Psi$.

238 Let T_k denote the time taken to compute task k , and $T_{k,j}$ represent the
239 time taken by UAV j to receive the data, process it and return the result
240 back to the master. We then have

$$T_k = \max_j T_{k,j} \quad (11)$$

241 For simplicity, we assume the result size is small and the latency of sending
242 it back is negligible, as often assumed in existing works [23]. Therefore, $T_{k,i}$
243 can be expressed as

$$T_{k,j} = T_{k,j}^{comp} + T_{k,0j}^{trans} \quad (12)$$

244 According to the computing and communication models described in the
245 previous section, we can derive that $T_{k,i}^{comp} = \frac{\xi_k c_i^k \ell_k}{f_i}$ and $T_{k,j}^{trans}$ satisfies

$$\begin{cases} \int_0^{T_{k,0j}^{trans}} \nu_{0j}(\tau) d\tau = c_j^k \ell_k & \text{if } j \neq 0 \\ T_{k,0j}^{trans} = 0 & \text{if } j = 0 \end{cases} \quad (13)$$

246 Concerning the energy consumption E_k of offloading task k , we focus on
247 the energy consumption of the entire NAC system, which consists of energy
248 consumed for computation and transmission among all the computing nodes
249 including both the master and workers, that is

$$E_k = \sum_{j=0}^N E_{k,j}^{comp} + \sum_{j=1}^N E_{k,0j}^{trans} \quad (14)$$

250 The Joint Optimization of Task Offloading and Resource Allocation for
251 Networked UAV Systems focuses on efficiently distributing computational

252 tasks and communication resources across a UAV network comprising a mas-
 253 ter UAV and multiple worker UAVs. In this system, the master UAV must
 254 decide how to partition its computational workload by determining the frac-
 255 tion of tasks to offload to each worker UAV and the fraction to process locally.
 256 At the same time, the master UAV must allocate available communication
 257 resources, including subcarriers (or channel bandwidth) and transmission
 258 power, while operating within the constraints of its energy and bandwidth
 259 budgets.

260 The objective of the joint optimization problem is defined to minimize
 261 the weighted sum of the system's total energy consumption and overall task
 262 completion latency, which can be mathematically formulated as follows

$$\underset{c_j^k, w_j^k, \psi_j^k}{\text{Minimize}} \quad \sum_{k=1}^K (1 - \eta)T_k + \eta E_k \quad (15a)$$

$$\text{subject to} \quad \sum_{j=0}^N c_j^k \ell_k = S_k, \quad \sum_{j=1}^N w_j^k = W, \quad \sum_{j=1}^N \psi_j^k = \Psi, \quad (15b)$$

$$c_j^k \in \mathbb{Z}, \quad 0 \leq c_j^k \leq L_k, \quad (15c)$$

$$w_j^k \in \mathbb{Z}, \quad 0 \leq w_j^k \leq W, \quad (15d)$$

$$\psi_j^k \in \mathbb{R}, \quad 0 \leq \psi_j^k \leq \Psi, \quad (15e)$$

$$\forall j \in \mathcal{N}, k \in \mathcal{K}. \quad (15f)$$

263 where $\eta \in [0, 1]$ is the weight parameter of time latency and energy con-
 264 sumption that control the trade-off between energy consumption and time
 265 latency.

266 4.2. Markov Decision Process

267 To solve the optimization problem formulated in the previous section, the
 268 greatest challenge lies in the unknown relationship between T_k , E_k and the
 269 decision variables c_j^k, w_j^k, ψ_j^k , as UAVs lack knowledge of the system models.
 270 Additionally, the randomness of UAVs' mobility presents another significant
 271 challenge.

272 The complexity of the optimization problem lies in the nonlinear and
 273 unknown relationship between the objective components T_k and E_k and the
 274 decision variables c_j^k, w_j^k , and ψ_j^k . Furthermore, the stochastic nature of UAV
 275 mobility and the lack of explicit system models make traditional optimization
 276 approaches unsuitable for solving this problem effectively. To address these

challenges, the optimization problem defined in Eq. 15 is reformulated as a
 277 Markov Decision Process (MDP), represented by the tuple $(\mathcal{S}, \mathcal{A}, r, P)$. This
 278 conversion provides a framework for modeling the sequential decision-making
 279 process in a dynamic environment. By employing reinforcement learning al-
 280 gorithms, particularly model-free approaches, the master UAV (as the agent)
 281 can iteratively learn an optimal policy. This eliminates the need for complete
 282 system knowledge and allows the agent to adapt to stochastic network dy-
 283 namics, enabling an effective solution for minimizing the weighted objective.
 284

285 *4.2.1. State*

286 We assume the master agent has access only to the real-time positions
 287 of all UAVs, denoted as $\mathbf{p}_j(t)$ for all $j \in \mathcal{N}$, at each discrete time step t .
 288 Additionally, the agent possesses prior knowledge of the observed idle com-
 289 putational resources $f_j(t)$ available on each UAV at time t . These inputs
 290 serve as the foundational information for dynamic task offloading and re-
 291 source allocation.

292 To capture the temporal dynamics of the UAV network, continuous time
 293 is discretized into intervals of fixed length Δt . Let t^k denote the start time
 294 of task k . At this point, the state is defined as a combination of the task's
 295 characteristics, the recent trajectories of all UAVs, and their respective con-
 296 figurations. Formally, the state at t^k is expressed as:

$$s_k = [S_k, (\boldsymbol{\rho}_0(k), f_0(t^k)), (\boldsymbol{\rho}_1(k), f_1(t^k)), \dots, (\boldsymbol{\rho}_N(k), f_N(t^k))] \in \mathcal{S},$$

297 where S_k represents the size of the incoming task k , $\boldsymbol{\rho}_j(k)$ denotes the recent
 298 trajectory of UAV j , and $f_j(t^k)$ corresponds to the observed idle computa-
 299 tional resources of UAV j at the start of task k .

300 The trajectory $\boldsymbol{\rho}_j(k)$ is constructed by stacking the $M + 1$ most recent
 301 positions of UAV j , including its position at t^k , to capture mobility patterns
 302 critical for estimating task execution and communication feasibility in the
 303 future. Mathematically, $\boldsymbol{\rho}_j(k)$ is defined as:

$$\boldsymbol{\rho}_j(k) = [\mathbf{p}_j(t^k - M\Delta t), \mathbf{p}_j(t^k - (M - 1)\Delta t), \dots, \mathbf{p}_j(t^k)].$$

304 Here, $\mathbf{p}_j(t)$ represents the coordinates of UAV j at time t .

305 This comprehensive state representation provides a detailed view of the
 306 system at time t^k , integrating spatial and computational aspects of the UAV
 307 network. The inclusion of task size S_k , UAV trajectories $\boldsymbol{\rho}_j(k)$, and idle com-
 308 putational resources $f_j(t^k)$ allows the master agent to make informed deci-

309 sions regarding task partitioning and resource allocation. By effectively cap-
 310 turing the pattern of system's dynamics, this state formulation ensures the
 311 agent can optimize task offloading policies in a highly dynamic and stochastic
 312 environment.

313 *4.2.2. Action*

314 Every time a task k arrives, the master UAV partitions it into sub-tasks
 315 and offloads them among all UAVs while allocating communication resources
 316 such as subcarriers and transmission power. The fraction of task k assigned
 317 to UAV $j \in \mathcal{N}$ is denoted by $\zeta_j^{k,d} \geq 0$, where the superscript d indicates
 318 data allocation. Similarly, the fraction of subcarriers allocated to UAV j
 319 for task k is represented as $\zeta_j^{k,s} \geq 0$, and the fraction of transmission power
 320 allocated to UAV j is denoted by $\zeta_j^{k,p} \geq 0$, where the superscripts s and
 321 p represent subcarrier and power allocation, respectively. These fractions
 322 satisfy the following constraints:

$$\sum_{j=0}^N \zeta_j^{k,d} = 1, \quad \sum_{j=1}^N \zeta_j^{k,s} = 1, \quad \text{and} \quad \sum_{j=1}^N \zeta_j^{k,p} = 1,$$

323 ensuring that the entire task, all subcarriers, and the total transmission power
 324 are fully allocated. The workload offloaded expressed as the number of atom
 325 tasks and subcarrier allocations expressed as the number of subcarriers are
 326 computed by rounding the fractional values. The workload offloaded to UAV
 327 j is given as:

$$c_j^k = \text{round} \left(\zeta_j^{k,d} L_k \right),$$

328 where S_k is the size of task k . The number of subcarriers allocated to UAV
 329 j is computed as

$$w_j^k = \text{round} \left(\zeta_j^{k,s} W \right),$$

330 where W is the total number of available subcarriers. The transmission power
 331 allocated to UAV j is retained as a fractional value:

$$\psi_j^k = \zeta_j^{k,p} \Psi,$$

332 where Ψ is the total transmission power budget. The action taken by the
 333 master UAV at the start time t^k is defined as:

$$a_k = \left[\zeta_0^{k,d}, \dots, \zeta_N^{k,d}; \zeta_1^{k,s}, \dots, \zeta_N^{k,s}; \zeta_1^{k,p}, \dots, \zeta_N^{k,p} \right] \in \mathcal{A},$$

334 where a combination of an N -simplex and two $(N-1)$ -simplices formed the
 335 action space \mathcal{A} .

336 4.2.3. Reward

337 To jointly minimize the total task completion time and energy consumption
 338 we define the reward function as a weighted sum of the improvement
 339 ratios in time latency and energy efficiency for each task k . The reward
 340 function is expressed as:

$$r(s_k, a_k) = (1 - \eta) \left(\frac{\tilde{T}_k - T_k}{\tilde{T}_k} \right) + \eta \left(\frac{\tilde{E}_k - E_k}{\tilde{E}_k} \right), \quad (16)$$

341 where $\eta \in [0, 1]$ is the weight parameter that controls the trade-off between
 342 energy consumption and time latency. A smaller η emphasizes minimizing
 343 task completion time, while a larger η prioritizes reducing energy consump-
 344 tion.

345 In this formulation, T_k represents the actual task completion time achieved
 346 through offloading, and $\tilde{T}_k = \frac{\xi_k S_k}{f_0}$ denotes the time required to execute task
 347 k locally at the master UAV. The term $\frac{\tilde{T}_k - T_k}{\tilde{T}_k}$ quantifies the relative improve-
 348 ment in task completion time achieved by offloading. Similarly, E_k is the
 349 energy consumed for task k with offloading, and \tilde{E}_k is the energy required
 350 to execute the task locally at the master UAV. The term $\frac{\tilde{E}_k - E_k}{\tilde{E}_k}$ quantifies
 351 the relative improvement in energy consumption. By combining these two
 352 ratios, the reward function ensures that both latency and energy efficiency
 353 are considered.

354 4.2.4. Transition

355 As discussed in Sec. 3, all UAVs move randomly according to predefined
 356 mobility models that incorporate collision avoidance schemes to ensure safe
 357 operation. The task at time t^k is completed within a total wall time T_k ,
 358 a random variable dependent on the current state s_k of all UAVs and the
 359 offloading action a_k . Consequently, the next state s_{k+1} , starting at $t^{k+1} =$
 360 $t^k + T_k$, is governed by the transition dynamics of the mobility model and
 361 the actions taken by the master UAV.

362 The transition model, which describes state transitions based on the cur-
 363 rent state and action, depends on the specific random mobility model. It is
 364 abstracted as:

$$s_{k+1} \sim \begin{cases} P(s_{k+1}|s_k, a_k; v_{\max}, \lambda_{\text{rd}}) & (\text{RD}) \\ P(s_{k+1}|s_k, a_k; \sigma, \lambda_{\text{st}}) & (\text{ST}) \end{cases} \quad (17)$$

365 Here, $P(\cdot)$ represents the transition probability, whose explicit form is
366 generally unknown and depends on the selected mobility model. For the
367 random direction (RD) model, v_{\max} controls the maximum velocity of the
368 UAVs, while λ_{rd} governs the frequency of direction changes. In the smooth
369 turn (ST) model, σ determines the radius of the circular trajectory, and λ_{st}
370 specifies the frequency of direction changes, thereby controlling how often
371 the UAV reorients its trajectory.

372 Although the explicit form of $P(\cdot)$ is generally unknown, it can be inferred
373 or approximated using observed UAV trajectories and environmental condi-
374 tions. This abstraction provides a robust framework for decision-making in
375 task allocation and resource management across dynamic, networked UAV
376 systems.

377 5. Deep Reinforcement Learning Solution

378 In this section, we introduce our model-free DRL-based algorithm, a vari-
379 ant of the Twin Delayed Deep Deterministic policy gradient algorithm (TD3)
380 [15], to address these challenges.

381 5.1. TD3-based Offloading

382 TD3 [15] is an advanced actor-critic algorithm designed for continuous
383 action spaces in DRL. Here, we introduce a variant of the TD3 algorithm to
384 enable the master UAV to identify trustworthy offloadees and optimize task
385 allocation.

386 5.1.1. Actor

387 The behavior of the master agent is defined by a policy function $\pi : \mathcal{S} \rightarrow$
388 \mathcal{A} , which maps each state $s \in \mathcal{S}$ to a continuous action $a \in \mathcal{A}$. The goal of the
389 agent is to determine the optimal policy π^* , which maximizes the expected
390 return defined as $J = \mathbb{E}_{s_k \sim P, a_k \sim \pi}[R_1]$, where $R_k = \sum_{i=k}^K \gamma^{i-k} r(s_i, a_i)$ is the
391 accumulative discounted reward and γ is the discount factor.

392 To approximate the policy function, TD3 [15] utilizes a neural network
393 parameterized by ϕ , denoted as π_ϕ , which directly maps the state space $s \in \mathcal{S}$
394 to the action space \mathcal{A} . To satisfy the constraints in (15b), we adjust the actor
395 π_ϕ as follows. By including all UAVs' historic trajectories in the state, the
396 neural network can learn the movement patterns of each UAV and their inter-
397 actions. Moreover, the network adjusts weights to prioritize UAVs that have
398 more resources to share and are more likely to remain close to the master.

399 Therefore, the neural network's output logits $\mathbf{z}^k \in \mathbb{R}^{N+1}$ can be interpreted
400 as the reliability score of each UAV. We then use the Softmax function to
401 decide the offloading portions according to UAVs' reliability scores as follows

$$\mu_i^k = \text{Softmax}(z_i^k) = \frac{e^{z_i^k}}{\sum_{j=0}^N e^{z_j^k}} \quad (18)$$

402 where z_i^k represents the reliability score of UAV i for completing task k .

403 To estimate the parameters ϕ , we apply off-policy learning to enhance
404 training stability. Moreover, to strengthen the robustness of the learned
405 policy function against variance and prevent overfitting to narrow peaks of
406 action values, we add random noise to the actions of the target actor $\pi_{\phi'}$
407 parameterized by ϕ' as follows [15]

$$\begin{aligned} \tilde{\mu}_i^k &= (1 - \beta)\mu_i^k + \beta\epsilon_i \\ \epsilon &\sim \text{Dir}(\alpha_{policy}) \end{aligned} \quad (19)$$

408 where β is the weight of noise ϵ_i , and $\epsilon_i \in \mathbb{R}^{N+1}$ is sampled from the Dirichlet
409 distribution [24] with a concentration parameter α_{policy} that is uniform across
410 all elements. Of note, the Dirichlet distribution guarantees that the actions
411 remain within the space of N -simplex.

412 5.1.2. Critic

413 Given the state s_k and the action a_k taken, the state-action value func-
414 tion Q^π provides the expected return when following policy π thereafter, i.e.,
415 $Q^\pi(s_k, a_k) = \mathbb{E}_{s_{i>k} \sim P, a_{i>k} \sim \pi}[R_k | s_k, a_k]$. To approximate the critic Q^π , neural
416 networks are also used. Following TD3[15], we define two primary critic net-
417 works Q_{θ_1} and Q_{θ_2} and two target critic networks $Q_{\theta'_1}$ and $Q_{\theta'_2}$ parameterized
418 by $\theta_1, \theta_2, \theta'_1$ and θ'_2 , respectively.

419 5.1.3. Training

420 As shown in Alg. 1, the training starts by initializing all the parameters
421 and the replay buffer \mathcal{D} , which stores transition samples (Lines 1-3). At
422 each training iteration u , the master agent interacts with the environment
423 to collect transition data and store them in the buffer \mathcal{D} until the number
424 of collected transitions exceeds H (Lines 5-6). After that, the agent utilizes
425 a batch \mathcal{B} sampled from the replay buffer to update the parameters of the
426 critics (Lines 9-11) and actor (Lines 12-13). Particularly, the parameters θ_i ,

⁴²⁷ $i \in \{1, 2\}$, of each primary critic is updated by minimizing the following loss
⁴²⁸ over the sampled batch

$$l_i = \frac{1}{|\mathcal{B}|} \sum_{(s,a,r,s') \in \mathcal{B}} (y - Q_{\theta_i}(s, a))^2 \quad (20)$$

⁴²⁹ where $y = r(s, a) + \gamma \min_{i=1,2} Q_{\theta'_i}(s', \tilde{a}')$ and $\tilde{a}' = (\tilde{\mu}_0, \tilde{\mu}_1, \dots, \tilde{\mu}_N)$ is the regu-
⁴³⁰ larized action defined in (19).

⁴³¹ A lower update frequency is necessary for the policy function compared to
⁴³² the value function, otherwise, it may lead to divergence. Hence, we update
⁴³³ the policy function every d iterations by maximizing the expected return
⁴³⁴ in the direction of the batch gradient of the policy, where the gradient is
⁴³⁵ computed by

$$\nabla_{\phi} J = \frac{1}{|\mathcal{B}|} \sum_{(s,a,r,s') \in \mathcal{B}} [\nabla_a Q_{\theta_1}(s, a)|_{a=\pi_{\phi}(s)} \nabla_{\phi} \pi_{\phi}(s)] \quad (21)$$

⁴³⁶ Also, the parameters of the target networks are gradually adjusted towards
⁴³⁷ the weights of the primary networks through weighted soft updates every d
⁴³⁸ iteration as follows

$$\begin{aligned} \theta'_i &\leftarrow \theta_i + (1 - \tau)(\theta'_i - \theta_i), i = 1, 2 \\ \phi' &\leftarrow \phi + (1 - \tau)(\phi' - \phi) \end{aligned} \quad (22)$$

⁴³⁹ where $\tau \in [0, 1]$ is the weight.

⁴⁴⁰ 6. Experiments

⁴⁴¹ In this section, we conduct experiments to evaluate the effectiveness and
⁴⁴² scalability of our proposed TD3-based DRL algorithm.

⁴⁴³ 6.1. Environment Setting

⁴⁴⁴ We evaluate our method in simulated scenarios with one master UAV
⁴⁴⁵ and $N = 3, 6, 9$, or 12 offload UAVs respectively. For the multi-UAV RD
⁴⁴⁶ mobility model, we set its parameters as $v_{max} = 20$ m/s, $\lambda = \frac{1}{15}$, and $D_i = 5$,
⁴⁴⁷ $\forall i \in \mathcal{N}$. The length of each discrete time step is set to $\Delta t = 1$ second. We
⁴⁴⁸ also vary the size of the flying zone and consider two sizes, $W = 300$ m and
⁴⁴⁹ $W = 400$ m. The computing power f_i of each UAV i is randomly configured
⁴⁵⁰ by selecting values from the range of [1, 1.6] Ghz. The task list consists of

Algorithm 1 TD3 training for task offloading

- 1: Initialize the critic networks $Q_{\theta_1}, Q_{\theta_2}$ by randomly assigning values to θ_1, θ_2 , and initialize the target critic networks $Q_{\theta'_1}, Q_{\theta'_2}$ by setting $\theta'_1 \leftarrow \theta_1, \theta'_2 \leftarrow \theta_2$
- 2: Initialize the actor network π_ϕ and the target actor $\pi_{\phi'}$ by randomly assigning values to ϕ and setting $\phi' \leftarrow \phi$
- 3: Initialize the replay buffer by $\mathcal{D} \leftarrow \emptyset$
- 4: **for** $u = 1$ to U **do**
- 5: Select action $a \sim (1 - \beta)\pi_\phi(s) + \beta\epsilon$, with exploration noise $\epsilon \sim \text{Dir}(\alpha_{explore})$, observe reward r and new state s'
- 6: Store the transition tuple (s, a, r, s') in \mathcal{D}
- 7: **if** $u > H$ **then**
- 8: Sample a batch of $|\mathcal{B}|$ transitions (s, a, r, s') from buffer \mathcal{D}
- 9: $\tilde{a}' \leftarrow (1 - \beta)\pi_{\phi'}(s') + \beta\epsilon$, $\epsilon \sim \text{Dir}(\alpha_{policy})$
- 10: $y \leftarrow r + \gamma \min_{i=1,2} Q_{\theta'_i}(s', \tilde{a}')$
- 11: Update the critics by $\theta_i \leftarrow \arg \min_{\theta_i} \frac{1}{|\mathcal{B}|} \sum (y - Q_{\theta_i}(s, a))^2$, $i = 1, 2$
- 12: **if** $u \bmod d = 0$ **then**
- 13: Update the actor by $\phi \leftarrow \arg \max_{\phi} \frac{1}{|\mathcal{B}|} \sum Q_{\theta_1}(s, \pi_\phi(s))$
- 14: Update the target network by $\theta'_i \leftarrow \theta_i + (1 - \tau)(\theta'_i - \theta_i)$, $i = 1, 2$
and $\phi' \leftarrow \phi + (1 - \tau)(\phi' - \phi)$
- 15: **end if**
- 16: **end if**
- 17: **end for**

451 $K = 25$ tasks that can be locally completed in $\tilde{T}_k \in [20, 60]$ seconds by the
 452 master UAV. All tasks can be divided into $L_k = 1000$ atomic tasks. The
 453 computation intensity is set to be the same $\xi_k = 10^6$ cycles/kB for all tasks,
 454 and the size of each task is determined by $S_k = \frac{f_0 \tilde{T}_k}{\xi_k}$. As a case of a networked
 455 UAV swarm utilizing 5G Wi-Fi communication for data transmission, we set
 456 the bandwidth $B = 40$ MHz, relative distance $d_r = 1$ meter, path gain
 457 $G = 40$, path loss exponent $\theta = 4$, and noise power spectral density $N_0 =$
 458 -174 dBm/Hz. The transmitted power ψ_{ij} of the master UAV to each UAV
 459 $i \in \mathcal{N} \setminus \{0\}$ varies from 80 mW to 120 mW.

460 6.2. Training Performance

461 We employ a 3-layer Multilayer Perceptrons (MLP)[25] architecture for
 462 both the actor and critic networks. The actor network takes observations as
 463 input, whose dimension is based on the number of computing nodes $N + 1$,
 464 and outputs actions of dimension N . Each UAV trajectory has a length
 465 of $M + 1$, where $M = 20$. The critic network takes the concatenation of
 466 observations and actions as input and outputs a scalar representing the state
 467 action value. The width of the hidden layer in both networks is set to twice
 468 the dimension of the input. All parameters are initialized using Kaiming
 469 initialization[26].

470 The learning rates for the actor and critic networks are set to 0.0001 and
 471 0.0002, respectively. The threshold H is set to 1000 and the batch size is
 472 $|\mathcal{B}| = 512$. The gradient norm is clipped between 0 and 0.2. The exploration
 473 and policy noise are sampled from a Dirichlet distribution with parameters
 474 $\alpha_{explore} = 0.1$ and $\alpha_{policy} = 0.99$, respectively. The noise weight is $\beta = 0.1$
 475 and the discount factor γ is 0.99. The actor network is updated every $d = 25$
 476 iterations, and the soft update weight τ for target networks is 0.005.

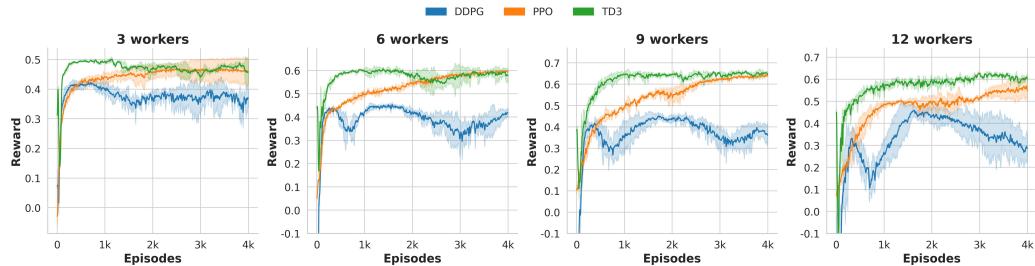


Figure 5: Learning Curve ($200 \times 200m^2$)

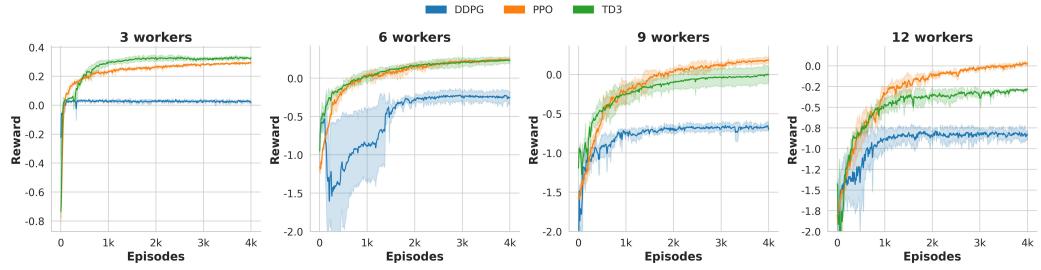


Figure 6: Learning Curve($300 \times 300m^2$)

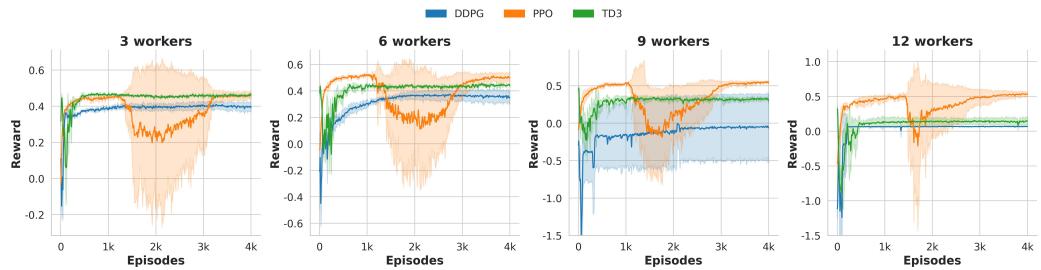


Figure 7: Learning Curve ($200 \times 200m^2$)

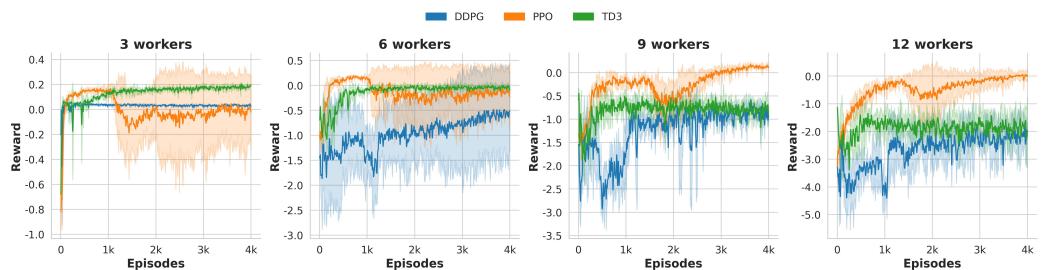


Figure 8: Learning Curve($300 \times 300m^2$)

In each scenario, we train the agent for 3000 episodes, i.e., $U = 3000K = 75000$ iterations, with three different random seeds. The results are shown in Fig. 7 and Fig. 8. The latency reduction is defined as the reduction in total task completion by task offloading compared to local computing at the master UAV, i.e.,

$$\text{Latency Reduction} = \frac{\sum_{k \in \mathcal{K}} \tilde{T}_k - \sum_{k \in \mathcal{K}} T_k}{\sum_{k \in \mathcal{K}} \tilde{T}_k} \times 100\%.$$

477 The two figures demonstrate that our method converges after training and
 478 shows promising stability. When UAVs are restricted to an area of $300 \times$
 479 $300m^2$ (Fig. 7), increasing the number of potential offloadees (workers) re-
 480 duces task completion time. The same phenomenon is observed when the
 481 flying zone expands to $400 \times 400m^2$, except when the number of potential
 482 workers increases to $N = 12$. The performance degradation at $N = 12$ may
 483 be due to increased collision avoidance maneuvers, which make UAVs' mobil-
 484 ity more uncertain and harder to learn and predict. Intuitively, the master
 485 agent tends to share workload with UAVs that are more likely to remain
 486 nearby throughout task execution, as indicated by a higher reliability score
 487 z_i^k . Therefore, when UAV mobility becomes more uncertain, fewer workers
 488 are selected as offloadees due to reduced reliability. This is confirmed by the
 489 results shown in Fig. 8, where performance improves when collision avoid-
 490 ance mechanisms are not in place. It also indicates that the task completion
 491 time cannot be infinitely reduced by continuously adding more UAVs to the
 492 region.

493 *6.3. Comparison Studies*

494 To evaluate the performance of the proposed method, we compare it
 495 with five benchmarks, including (1) **Equal (all)** that equally divides and
 496 distributes tasks to all UAVs; (2) **Equal (close)** that equally partitions
 497 and distributes tasks to UAVs that are close enough with distance to the
 498 master satisfying $d_i < 100m$; (3) **Naive Prediction** that selects offloadees
 499 based on predicted future trajectories and assigns sub-tasks of equal size to
 500 these offloadees. Specifically, it predicts each UAV's trajectory for the next
 501 20 steps, assuming the UAVs maintain their current velocity and ignoring
 502 potential collisions. It then calculates the percentage q_i of predicted UAV
 503 positions that remain within 200m of the master ($d_i < 200m$). UAVs with
 504 $q_i \geq 40\%$ are selected as offloadees; (4) **Reliable** that allocates tasks based

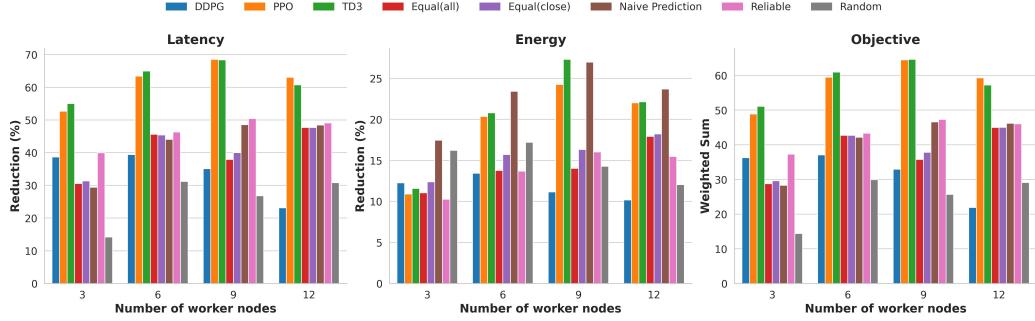


Figure 9: Performance Comparison ($200 \times 200m^2$)

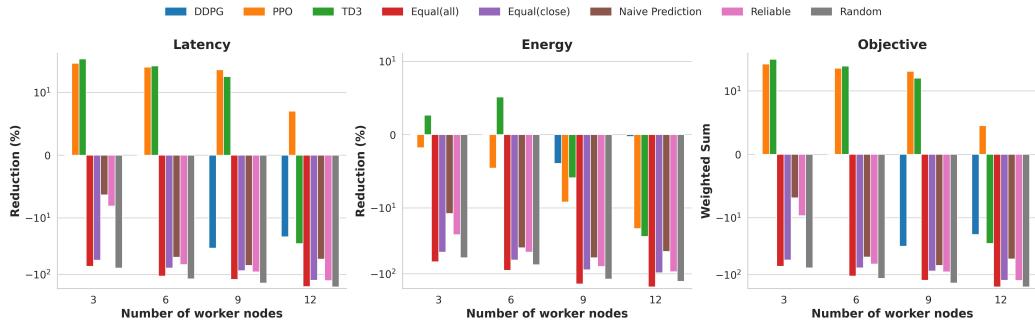


Figure 10: Performance Comparison ($300 \times 300m^2$)

on a reliability score defined as $\text{reliability} = q_i \psi_{ij} f_i$. It picks the top $\lceil \frac{N+1}{2} \rceil$ UAVs with the highest reliability scores and assigns tasks to these UAVs proportionally to their reliability scores; (5) **Random** that randomly sample actions from the action space.

The results in Fig. 11 and Fig. 12 demonstrate the promising performance of our method, which achieves the highest latency reduction across all scenarios. When the area is $300 \times 300m^2$, the **Naive Prediction** ranks second in scenarios with 3, 6, and 9 workers, while the **Equal (close)** ranks second in scenarios with 12 workers. When the area expands to $400 \times 400m^2$ (Fig. 12), the **Equal (all)** and the **Random** provide no benefit in reducing latency; instead, they significantly delay task completion. Comparing Fig. 11 and Fig. 12, we can observe that as the area increases for a fixed N , the performance of all methods degrades. This is due to the sparser airspace, which causes UAVs to be farther apart from each other, thereby increasing

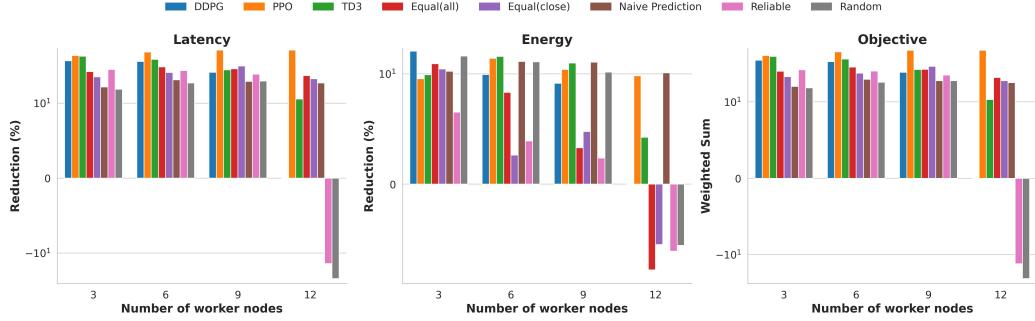


Figure 11: Performance Comparison ($200 \times 200m^2$)

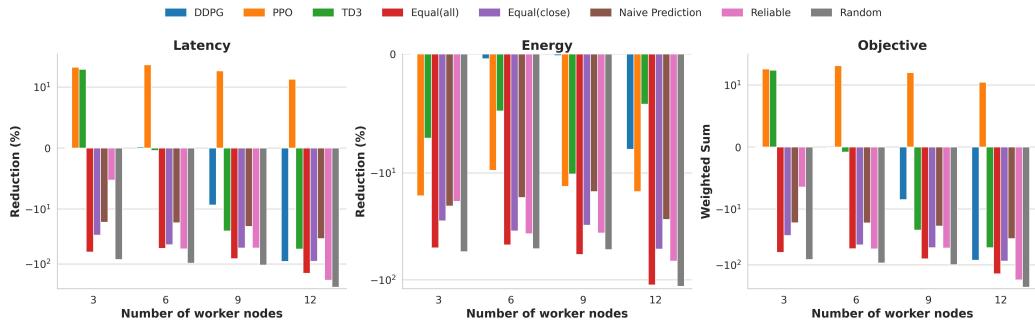


Figure 12: Performance Comparison ($300 \times 300m^2$)

519 transmission latency and task completion time.

520 Appendix A. Example Appendix Section

521 Appendix text.

522 Example citation, See [4].

523 References

- 524 [1] H. Kurunathan, H. Huang, K. Li, W. Ni, E. Hossain, Machine learning-
525 aided operations and communications of unmanned aerial vehicles:
526 A contemporary survey, IEEE Communications Surveys & Tutorials
527 (2023).

- 528 [2] Y. Zeng, R. Zhang, T. J. Lim, Wireless communications with unmanned
529 aerial vehicles: Opportunities and challenges, *IEEE Communications*
530 magazine 54 (5) (2016) 36–42.
- 531 [3] B. Liu, W. Zhang, W. Chen, H. Huang, S. Guo, Online computation
532 offloading and traffic routing for uav swarms in edge-cloud computing,
533 *IEEE Transactions on Vehicular Technology* 69 (8) (2020) 8777–8791.
- 534 [4] Z. Bai, Y. Lin, Y. Cao, W. Wang, Delay-aware cooperative task offload-
535 ing for multi-uav enabled edge-cloud computing, *IEEE Transactions on*
536 *Mobile Computing* (2022).
- 537 [5] M. Satyanarayanan, The emergence of edge computing, *Computer* 50 (1)
538 (2017) 30–39.
- 539 [6] Q. Hu, Y. Cai, G. Yu, Z. Qin, M. Zhao, G. Y. Li, Joint offloading and
540 trajectory design for uav-enabled mobile edge computing systems, *IEEE*
541 *Internet of Things Journal* 6 (2) (2018) 1879–1892.
- 542 [7] Y. Miao, K. Hwang, D. Wu, Y. Hao, M. Chen, Drone swarm path plan-
543 ning for mobile edge computing in industrial internet of things, *IEEE*
544 *Transactions on Industrial Informatics* (2022).
- 545 [8] K. Lu, J. Xie, Y. Wan, S. Fu, Toward uav-based airborne computing,
546 *IEEE Wireless Communications* 26 (6) (2019) 172–179.
- 547 [9] H. Zhang, B. Wang, R. Wu, J. Xie, Y. Wan, S. Fu, K. Lu, Exploring net-
548 worked airborne computing: A comprehensive approach with advanced
549 simulator and hardware testbed, *Unmanned Systems* (2023).
- 550 [10] B. Wang, J. Xie, K. Lu, Y. Wan, S. Fu, Learning and batch-processing
551 based coded computation with mobility awareness for networked air-
552 borne computing, *IEEE Transactions on Vehicular Technology* (2022).
- 553 [11] E. Shtaiwi, A. Abdelhadi, H. Li, Z. Han, H. V. Poor, Orthogonal time
554 frequency space for integrated sensing and communication: A survey,
555 arXiv preprint arXiv:2402.09637 (2024).
- 556 [12] W. Lu, P. Si, Y. Gao, H. Han, Z. Liu, Y. Wu, Y. Gong, Trajectory and
557 resource optimization in ofdm-based uav-powered iot network, *IEEE*
558 *Transactions on Green Communications and Networking* 5 (3) (2021)
559 1259–1270.

- 560 [13] X. Guan, Y. Huang, Q. Shi, Joint subcarrier and power allocation for
561 multi-uav systems, *China Communications* 16 (1) (2019) 47–56.
- 562 [14] E. M. Royer, P. M. Melliar-Smith, L. E. Moser, An analysis of the
563 optimum node density for ad hoc mobile networks, in: *ICC 2001. IEEE*
564 *International Conference on Communications. Conference Record (Cat.*
565 *No. 01CH37240)*, Vol. 3, IEEE, 2001, pp. 857–861.
- 566 [15] S. Fujimoto, H. van Hoof, D. Meger, Addressing Function Approxima-
567 *tion Error in Actor-Critic Methods* (Oct. 2018). arXiv:1802.09477.
- 568 [16] K. Cheng, Y. Teng, W. Sun, A. Liu, X. Wang, Energy-efficient joint
569 offloading and wireless resource allocation strategy in multi-mec server
570 systems, in: *2018 IEEE international conference on communications* (ICC),
571 IEEE, 2018, pp. 1–6.
- 572 [17] F. Pervez, A. Sultana, C. Yang, L. Zhao, Energy and latency efficient
573 joint communication and computation optimization in a multi-uav as-
574 sisted mec network, *IEEE Transactions on Wireless Communications*
575 (2023).
- 576 [18] A. Goldsmith, *Wireless communications*, Cambridge university press,
577 2005.
- 578 [19] X. Zhang, J. Xie, Drl-based task offloading for networked uavs with
579 random mobility and collision avoidance, in: *2024 20th International*
580 *Conference on Wireless and Mobile Computing, Networking and Com-*
581 *munications (WiMob)*, IEEE, 2024, pp. 514–519.
- 582 [20] D. H. Choi, S. H. Kim, D. K. Sung, Energy-efficient maneuvering
583 and communication of a single uav-based relay, *IEEE Transactions on*
584 *Aerospace and Electronic Systems* 50 (3) (2014) 2320–2327.
- 585 [21] Y. Wan, K. Namuduri, Y. Zhou, D. He, S. Fu, A smooth-turn mobility
586 model for airborne networks, in: *Proceedings of the first ACM MobiHoc*
587 *workshop on Airborne Networks and Communications*, 2012, pp. 25–30.
- 588 [22] D. S. Lakew, U. Sa’ad, N.-N. Dao, W. Na, S. Cho, Routing in flying ad
589 *hoc networks: A comprehensive survey*, *IEEE Communications Surveys*
590 & *Tutorials* 22 (2) (2020) 1071–1120.

- 591 [23] N. T. Hoa, N. C. Luong, D. Van Le, D. Niyato, et al., Deep reinforcement
592 learning for multi-hop offloading in uav-assisted edge computing, IEEE
593 Transactions on Vehicular Technology (2023).
- 594 [24] C. M. Bishop, N. M. Nasrabadi, Pattern recognition and machine learn-
595 ing, Vol. 4, Springer, 2006.
- 596 [25] D. E. Rumelhart, G. E. Hinton, R. J. Williams, Learning representations
597 by back-propagating errors, nature 323 (6088) (1986) 533–536.
- 598 [26] K. He, X. Zhang, S. Ren, J. Sun, Delving deep into rectifiers: Surpassing
599 human-level performance on imagenet classification, in: Proceedings of
600 the IEEE international conference on computer vision, 2015, pp. 1026–
601 1034.