

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-

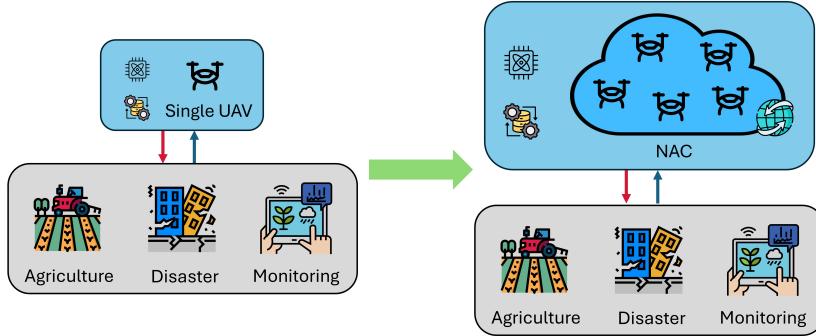


Figure 1: Networked Airborne Computing (NAC)

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. 2 details the system models and formulates
95 the optimization problem. Sec. 4 describes the proposed DRL algorithm.

⁹⁶ In Sec. 5, simulation results are presented and discussed. We conclude in
⁹⁷ Sec. ??.

⁹⁸ 2. System Models

⁹⁹ In this section, we will introduce the system models. Consider a group of
¹⁰⁰ $N + 1$ heterogeneous UAVs with varying physical configurations, indexed as
¹⁰¹ $i \in \mathcal{N} = \{0, 1, 2, \dots, N\}$. Each UAV is equipped with computing and commu-
¹⁰² nication modules, enabling resource sharing and onboard computation. Their
¹⁰³ characteristics of computing, communication, energy consumption, and mo-
¹⁰⁴ bility within the system can be comprehensively described and modeled as
¹⁰⁵ follows

¹⁰⁶ 2.1. Computing Model

¹⁰⁷ We describe the computing capability of each UAV i as CPU cycle fre-
¹⁰⁸ quency f_i in Hz. For a general computing task k , its input data size is S_k
¹⁰⁹ (bits), and its required computation intensity is ξ_k (cycles/bit)[16]. The total
¹¹⁰ CPU cycles required to compute task k is hence $\xi_k S_k$ and the time required
¹¹¹ for UAV i to execute this task is

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

¹¹² The corresponding energy consumption during computing can be given as

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

¹¹³ where ϵ represents an energy consumption parameter associated with the
¹¹⁴ effective switched capacitance, which is determined by the underlying CPU
¹¹⁵ architecture. To simplify the analysis, it is assumed that this capacitance
¹¹⁶ remains uniform across all devices [17].

¹¹⁷ 2.2. Communication Model

¹¹⁸ Let the distance between UAV i and UAV j be denoted as d_{ij} . The
¹¹⁹ UAV-to-UAV links are typically Line of Sight (LoS), with propagation speed
¹²⁰ approaching the speed of light. Hence, the transmission latency can be ap-
¹²¹ proximated using the transmission time. Here, we model the transmission
¹²² rate (bits/s) based on the Simplified Path Loss Model [18] as follows

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

123 where B_{ij} is the bandwidth (Hz) of the specific channel between UAV i
 124 and j , d_r is constant reference distance (meter), G is the unitless constant
 125 equal to the path gain of the distance d_r , θ is the path loss exponent, ψ_i is
 126 the transmitted power (mW) and N_0 is the constant noise power spectral
 127 density (dBm/Hz), ν_{ij} is the threshold to suppress the data rate too low,
 128 regarded as a constraint on quality of communication.

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

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

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

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

137 2.3. Mobility Model

138 We assume the UAVs fly at the same altitude. Therefore, the position
 139 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
 140 constraints $0 \leq x_i(t) \leq W$, $0 \leq y_i(t) \leq W$, such that the position is bounded
 141 within an area of $W \times W(m^2)$. Given initial position $\mathbf{p}_i(0) = (x_i(0), y_i(0))$,
 142 we adapt the Random Direction (RD) model [14] and Smooth Turn (ST)
 143 model[21] to model its movement, which have been widely used for describing
 144 UAVs, particularly multirotor drones [22]. At time t , let the velocity of UAV
 145 i be $\mathbf{v}_i(t) = (v_{i,x}(t), v_{i,y}(t)) \in \mathbb{R}^2$ and the heading direction be $\Phi_i(t) \in [0, 2\pi]$,
 146 where $v_{i,x}(t)$, $v_{i,y}(t)$ are component of velocity in x and y direction.

147 2.3.1. Random Direction

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

152 Suppose the start time of m -th duration is denoted as $t_m = \sum_{l=0}^{m-1} \delta_{i,l}$
 153 which is also the m -th instance when UAV i changes its velocity, where each

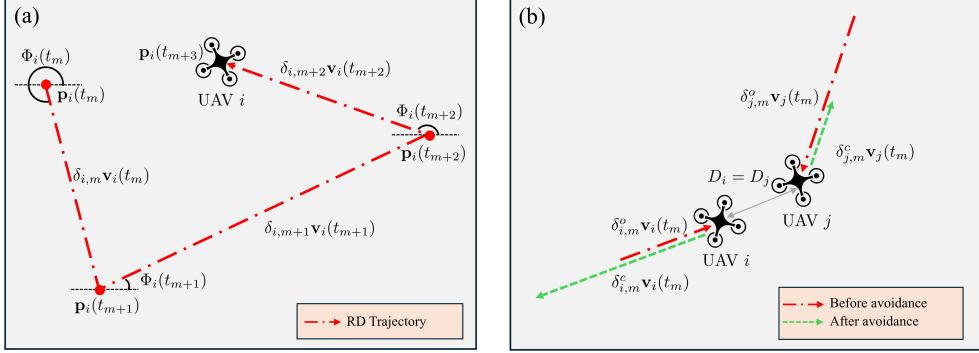


Figure 2: Caption

154 $\delta_{i,l}$ is randomly sampled from exponential distribution with parameter λ_{rd}
 155 and t_0 is set to be 0. At time t_m , UAV i select a speed magnitude uniformly
 156 between 0 and $v_{max} \in \mathbb{R}$, i.e. $0 < |\mathbf{v}_i(t_m)| < v_{max}$ and the heading direction
 157 $\Phi_i(t_m)$ in radian uniformly across 2π , leading to the new velocity $\mathbf{v}_i(t_m)$ for
 158 the m -th duration such that

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

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

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

161 The traditional RD model [14] is originally designed to describe the mo-
 162 bility of independent single entities, which ignores the spatiotemporal cor-
 163 relations of trajectory across entities. However, in multi-UAV systems, the
 164 mobility of UAVs can change to avoid collisions. This necessitates the in-
 165 corporation of collision avoidance mechanisms into the RD model. Here, we
 166 assume that each UAV i will turn around and move in the opposite direction
 167 without changing speed until the current duration is completed when its dis-
 168 tance to any other UAV j falls below a threshold D_i . Note that if both UAVs

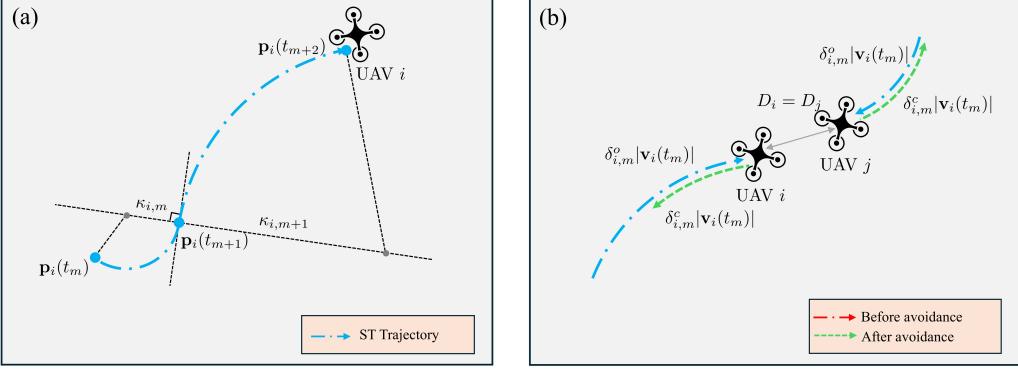


Figure 3: Caption

have the same threshold $D_i = D_j$, UAV j will also reserve its direction. By treating the boundaries of the area as obstacles, we ensure that all UAVs move within the designated area. The mobility of each UAV i with collision avoidance can then be described as

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

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

2.3.2. Smooth Turn

In the smooth turn (ST) mobility model, UAV i randomly picks a turning center located at a point along the line perpendicular to its current heading direction and moves at a constant forward speed following the circular trajectory around the turning center for a duration randomly selected until another turning center picked and repeat the process.

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

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

and the velocity follows as

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

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

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

The same issues appear in the original ST mobility model as the RD mobility model since it did not include collision avoidance mechanisms either. Here, we propose a collision avoidance strategy similar to the one introduced for the RD model in the previous section. In this approach, UAVs retrace their prior circular trajectory in the reverse direction, restoring the safe inter-UAV spacing and spatial relationship. The heading angle $\Phi_i(t)$ in Eq. 8 of UAV i will become

$$\Phi_i(t) = \begin{cases} \Phi_i(t_m) - \frac{|\mathbf{v}_i(t)|}{\kappa_{i,m}}(t - t_m) \\ \Phi_i(t_m) + \frac{|\mathbf{v}_i(t)|}{\kappa_{i,m}}(t - t_m^c) + \pi \end{cases} \quad (10)$$

where the superscript "o" represents the original value and the superscript "c" represents the value changed after the collision avoidance maneuver. Then the mobility of each UAV i equipped with a collision avoidance mechanism can modeled as

207 **3. Problem Formulation**

208 Without loss of generality, we let UAV $i = 0$ be the master (or offloader)
 209 and treat the remaining N UAVs as potential offloadees with idle computing
 210 resources. Suppose a sequence of computing tasks $\mathcal{K} = \{1, 2, \dots, K\}$ is gener-
 211 ated at the master, and each task k can be divided into $L_k \in \mathbb{Z}^+$ atomic tasks
 212 of size $\ell_k = \frac{S_k}{L_k}$, which can be computed in parallel. To minimize the total
 213 task completion time, the master aims to optimally allocate the L_k atomic
 214 tasks among all the UAVs, including itself. Suppose the workload offloaded
 215 to UAV $i \in \mathcal{N}$ for task k is $c_i^k \ell_k$, where c_i^k is a non-negative integer that
 216 satisfies $0 \leq c_i^k \leq L_k$. Therefore, $\sum_{i=0}^N c_i^k = L_k$. Of note, when there is no
 217 workload offloaded to UAV i , then $c_i^k = 0$.

218 Let T_k denote the time taken to compute task k , and $T_{k,i}$ represent the
 219 time taken by UAV i to receive the data, process it and return the result
 220 back to the master. We then have

$$T_k = \max_i T_{k,i} \quad (11)$$

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

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

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

$$\begin{cases} \int_0^{T_{k,i}^{trans}} \nu_{0i}(t) dt = c_i^k \ell_k & \text{if } i \neq 0 \\ T_{k,i}^{trans} = 0 & \text{if } i = 0 \end{cases} \quad (13)$$

226 The problem can then be mathematically formulated as follows

$$\underset{c_i^k, \forall i \in \mathcal{N}, k \in \mathcal{K}}{\text{Minimize}} \quad \sum_{k=1}^K T_k \quad (14a)$$

$$\text{subject to} \quad \sum_{i=0}^N c_i^k \ell_k = S_k \quad \forall k \in \mathcal{K} \quad (14b)$$

$$c_i^k \in \mathbb{Z}, 0 \leq c_i^k \leq L_k, \forall i \in \mathcal{N}, k \in \mathcal{K} \quad (14c)$$

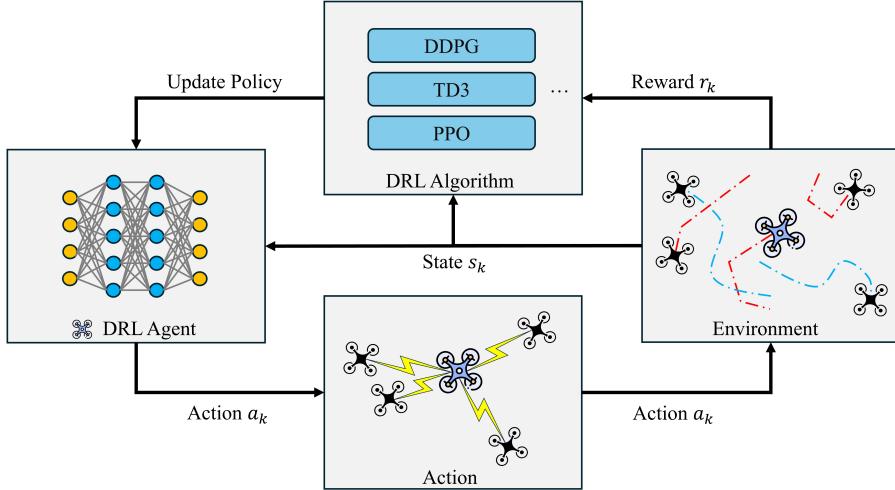


Figure 4: Caption

227 4. Deep Reinforcement Learning Solution

228 To solve the optimization problem formulated in the previous section, the
 229 greatest challenge lies in the unknown relationship between T_k and the decision
 230 variables c_i^k , as UAVs lack knowledge of the system models. Additionally,
 231 the randomness of UAVs' mobility presents another significant challenge. In
 232 this section, we introduce our model-free DRL-based algorithm, a variant of
 233 the Twin Delayed Deep Deterministic policy gradient algorithm (TD3) [15],
 234 to address these challenges.

235 *4.1. Markov Decision Process*

236 We first convert the optimization problem into a Markov Decision Process
 237 represented by a tuple $(\mathcal{S}, \mathcal{A}, r, P)$, where \mathcal{S} is the state space, \mathcal{A} is the
 238 action space, P is the state transition model, and r is the reward function.
 239 The master UAV is the DRL agent that takes actions to minimize the total
 240 task completion time.

241 *4.1.1. State*

242 We assume the master agent only has access to all UAVs' positions
 243 $\mathbf{p}_i(t), \forall i \in \mathcal{N}$ at each time t , and the prior knowledge about their config-
 244 urations, *i.e.* idle computing resources f_i , and communication capabilities
 245 characterized by transmitted powers ψ_i . We discretize the time into steps of

length Δt . Let t^k denote the start time of task k , the state at t^k is defined as the combination of the task size S_k , all UAVs' historic trajectories and configurations $s_k = [S_k, (\rho_0(k), f_0, \psi_0), (\rho_1(k), f_1, \psi_1), \dots, (\rho_N(k), f_N, \psi_N)] \in \mathcal{S}$, where $\rho_i(k) = [\mathbf{p}_i(t^k - M\Delta t), \mathbf{p}_i(t^k - (M-1)\Delta t), \dots, \mathbf{p}_i(t^k)] \in \mathbb{R}^{2 \times (M+1)}$ is UAV i 's trajectory consisting of its most recent M -step positions (including the start point).

4.1.2. Action

Every time a task k arrives, the master agent partitions it into sub-tasks of varying amounts of atomic tasks and offloads them to different UAVs. Let $\mu_i^k \geq 0$ denote the portion of task k offloaded to UAV $i \in \mathcal{N}$, such that $\sum_{i=0}^N \mu_i^k = 1$. The action taken by the master at time t^k is defined as $a_k = [\mu_0^k, \mu_1^k, \dots, \mu_N^k] \in \mathcal{A}$, where \mathcal{A} is the space of N -simplex. The workload offloaded to UAV i is then computed by $c_i^k = \text{round}(\frac{\mu_i^k S_k}{\ell_k})$.

4.1.3. Reward

To minimize the total task completion time, we define the reward function as follows

$$r(s_k, a_k) = \frac{S_k}{\sum_{k \in \mathcal{K}} S_k} \cdot \frac{\tilde{T}_k - T_k}{\tilde{T}_k} \quad (15)$$

where $\tilde{T}_k = \frac{\xi_k S_k}{f_0}$ represents the time required to execute task k locally at the master and $\frac{T_k - \tilde{T}_k}{T_k}$ indicates the acceleration rate achieved by offloading the task to nearby UAVs. $\frac{S_k}{\sum_{k \in \mathcal{K}} S_k}$ is the weight of task k among all tasks.

4.1.4. Transition

As discussed in Sec. 2, all UAVs move randomly according to the random mobility model with collision avoidance schemes. Hence, the transition model, which describes the state transitions given the current action, depends on the random mobility model and can be abstracted as follows

$$s_{k+1} \sim P(s_{k+1} | s_k, a_k; v_{max}, \lambda) \quad (16)$$

where $P(\cdot)$ represents the transition model, whose explicit form is unknown. v_{max} and λ are parameters of the random mobility model.

272 4.2. TD3-based Offloading

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

277 4.2.1. Actor

278 The behavior of the master agent is defined by a policy function $\pi : \mathcal{S} \rightarrow$
 279 \mathcal{A} , which maps each state $s \in \mathcal{S}$ to a continuous action $a \in \mathcal{A}$. The goal of the
 280 agent is to determine the optimal policy π^* , which maximizes the expected
 281 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
 282 accumulative discounted reward and γ is the discount factor.

283 To approximate the policy function, TD3 [15] utilizes a neural network
 284 parameterized by ϕ , denoted as π_ϕ , which directly maps the state space $s \in \mathcal{S}$
 285 to the action space \mathcal{A} . To satisfy the constraints in (14b), we adjust the actor
 286 π_ϕ as follows. By including all UAVs' historic trajectories in the state, the
 287 neural network can learn the movement patterns of each UAV and their inter-
 288 actions. Moreover, the network adjusts weights to prioritize UAVs that have
 289 more resources to share and are more likely to remain close to the master.
 290 Therefore, the neural network's output logits $\mathbf{z}^k \in \mathbb{R}^{N+1}$ can be interpreted
 291 as the reliability score of each UAV. We then use the Softmax function to
 292 decide the offloading portions according to UAVs' reliability scores as follows

$$\mu_i^k = \text{Softmax}(z_i^k) = \frac{e^{z_i^k}}{\sum_{j=0}^N e^{z_j^k}} \quad (17)$$

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

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

$$\begin{aligned} \tilde{\mu}_i^k &= (1 - \beta)\mu_i^k + \beta\epsilon_i \\ \epsilon &\sim \text{Dir}(\alpha_{policy}) \end{aligned} \quad (18)$$

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

303 4.2.2. Critic

304 Given the state s_k and the action a_k taken, the state-action value function
 305 Q^π provides the expected return when following policy π thereafter, i.e.,
 306 $Q^\pi(s_k, a_k) = \mathbb{E}_{s_{i>k} \sim P, a_{i>k} \sim \pi}[R_k | s_k, a_k]$. To approximate the critic Q^π , neural
 307 networks are also used. Following TD3[15], we define two primary critic net-
 308 works Q_{θ_1} and Q_{θ_2} and two target critic networks $Q_{\theta'_1}$ and $Q_{\theta'_2}$ parameterized
 309 by $\theta_1, \theta_2, \theta'_1$ and θ'_2 , respectively.

310 4.2.3. Training

311 As shown in Alg. 1, the training starts by initializing all the parameters
 312 and the replay buffer \mathcal{D} , which stores transition samples (Lines 1-3). At
 313 each training iteration u , the master agent interacts with the environment
 314 to collect transition data and store them in the buffer \mathcal{D} until the number
 315 of collected transitions exceeds H (Lines 5-6). After that, the agent utilizes
 316 a batch \mathcal{B} sampled from the replay buffer to update the parameters of the
 317 critics (Lines 9-11) and actor (Lines 12-13). Particularly, the parameters θ_i ,
 318 $i \in \{1, 2\}$, of each primary critic is updated by minimizing the following loss
 319 over the sampled batch

$$l_i = \frac{1}{|\mathcal{B}|} \sum_{(s, a, r, s') \in \mathcal{B}} (y - Q_{\theta_i}(s, a))^2 \quad (19)$$

320 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-
 321 larized action defined in (18).

322 A lower update frequency is necessary for the policy function compared to
 323 the value function, otherwise, it may lead to divergence. Hence, we update
 324 the policy function every d iterations by maximizing the expected return
 325 in the direction of the batch gradient of the policy, where the gradient is
 326 computed by

$$\nabla_\phi J = \frac{1}{|\mathcal{B}|} \sum_{(s, a, r, s') \in \mathcal{B}} [\nabla_a Q_{\theta_1}(s, a)|_{a=\pi_\phi(s)} \nabla_\phi \pi_\phi(s)] \quad (20)$$

327 Also, the parameters of the target networks are gradually adjusted towards
 328 the weights of the primary networks through weighted soft updates every d
 329 iteration as follows

$$\begin{aligned} \theta'_i &\leftarrow \theta_i + (1 - \tau)(\theta'_i - \theta_i), i = 1, 2 \\ \phi' &\leftarrow \phi + (1 - \tau)(\phi' - \phi) \end{aligned} \quad (21)$$

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

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**

331 **5. Experiments**

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

334 *5.1. Environment Setting*

335 We evaluate our method in simulated scenarios with one master UAV
336 and $N = 3, 6, 9$, or 12 offload UAVs respectively. For the multi-UAV RD
337 mobility model, we set its parameters as $v_{max} = 20$ m/s, $\lambda = \frac{1}{15}$, and $D_i = 5$,
338 $\forall i \in \mathcal{N}$. The length of each discrete time step is set to $\Delta t = 1$ second. We
339 also vary the size of the flying zone and consider two sizes, $W = 300$ m and
340 $W = 400$ m. The computing power f_i of each UAV i is randomly configured
341 by selecting values from the range of $[1, 1.6]$ Ghz. The task list consists of
342 $K = 25$ tasks that can be locally completed in $\tilde{T}_k \in [20, 60]$ seconds by the
343 master UAV. All tasks can be divided into $L_k = 1000$ atomic tasks. The
344 computation intensity is set to be the same $\xi_k = 10^6$ cycles/kB for all tasks,
345 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
346 UAV swarm utilizing 5G Wi-Fi communication for data transmission, we set
347 the bandwidth $B = 40$ MHz, relative distance $d_r = 1$ meter, path gain
348 $G = 40$, path loss exponent $\theta = 4$, and noise power spectral density $N_0 =$
349 -174 dBm/Hz. The transmitted power ψ_i of the master UAV to each UAV
350 $i \in \mathcal{N} \setminus \{0\}$ varies from 80 mW to 120 mW.

351 *5.2. Training Performance*

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

361 The learning rates for the actor and critic networks are set to 0.0001 and
362 0.0002 , respectively. The threshold H is set to 1000 and the batch size is
363 $|\mathcal{B}| = 512$. The gradient norm is clipped between 0 and 0.2 . The exploration
364 and policy noise are sampled from a Dirichlet distribution with parameters
365 $\alpha_{explore} = 0.1$ and $\alpha_{policy} = 0.99$, respectively. The noise weight is $\beta = 0.1$

366 and the discount factor γ is 0.99. The actor network is updated every $d = 25$
367 iterations, and the soft update weight τ for target networks is 0.005.

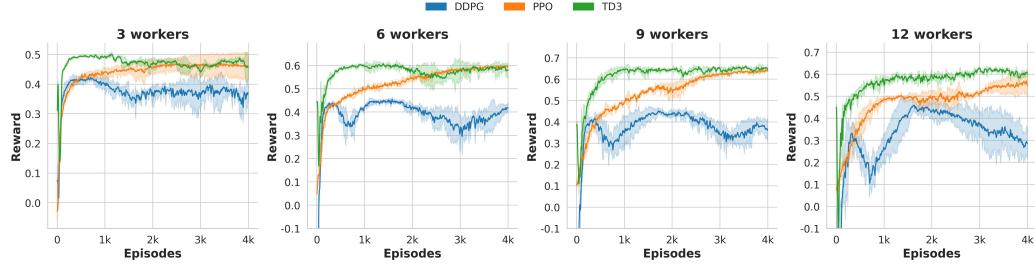


Figure 5: Learning Curve ($200 \times 200m^2$)

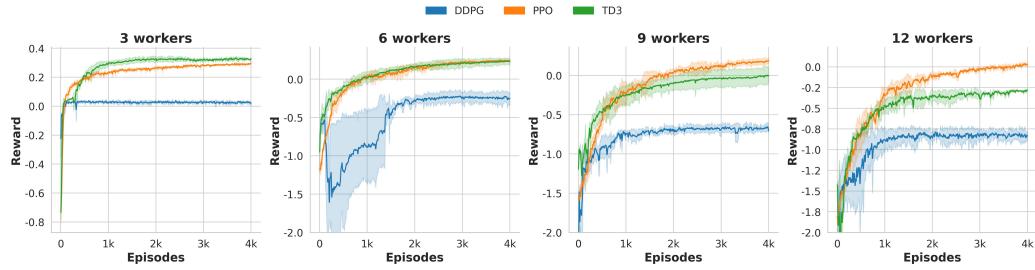


Figure 6: Learning Curve($300 \times 300m^2$)

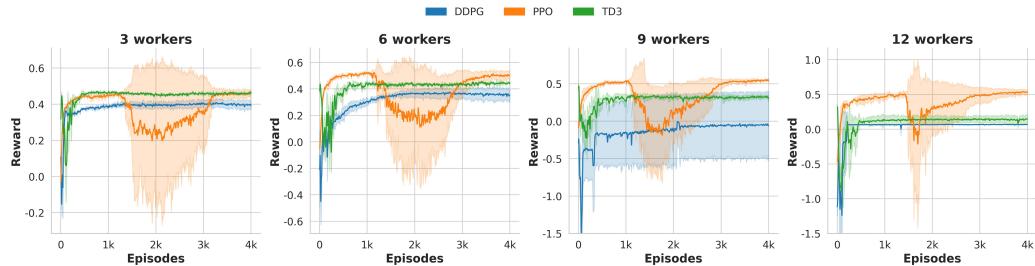


Figure 7: Learning Curve ($200 \times 200m^2$)

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

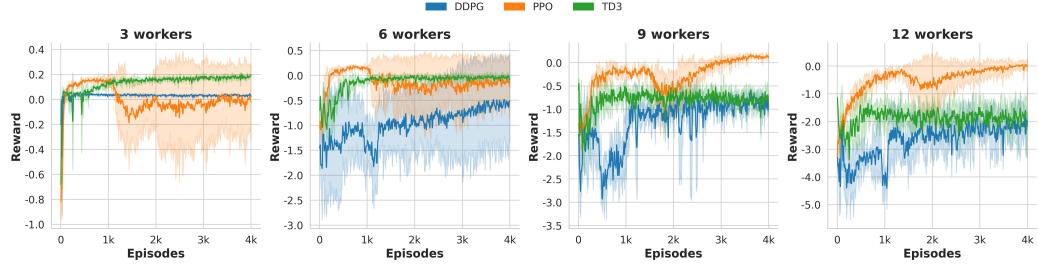


Figure 8: Learning Curve($300 \times 300m^2$)

total task completion by task offloading compared to local computing at the master UAV, i.e.,

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

368 The two figures demonstrate that our method converges after training and
 369 shows promising stability. When UAVs are restricted to an area of $300 \times$
 370 $300m^2$ (Fig. 7), increasing the number of potential offloadees (workers) re-
 371 duces task completion time. The same phenomenon is observed when the
 372 flying zone expands to $400 \times 400m^2$, except when the number of potential
 373 workers increases to $N = 12$. The performance degradation at $N = 12$ may
 374 be due to increased collision avoidance maneuvers, which make UAVs' mobil-
 375 ity more uncertain and harder to learn and predict. Intuitively, the master
 376 agent tends to share workload with UAVs that are more likely to remain
 377 nearby throughout task execution, as indicated by a higher reliability score
 378 z_i^k . Therefore, when UAV mobility becomes more uncertain, fewer workers
 379 are selected as offloadees due to reduced reliability. This is confirmed by the
 380 results shown in Fig. 8, where performance improves when collision avoid-
 381 ance mechanisms are not in place. It also indicates that the task completion
 382 time cannot be infinitely reduced by continuously adding more UAVs to the
 383 region.

384 *5.3. Comparison Studies*

385 To evaluate the performance of the proposed method, we compare it
 386 with five benchmarks, including (1) **Equal (all)** that equally divides and
 387 distributes tasks to all UAVs; (2) **Equal (close)** that equally partitions
 388 and distributes tasks to UAVs that are close enough with distance to the

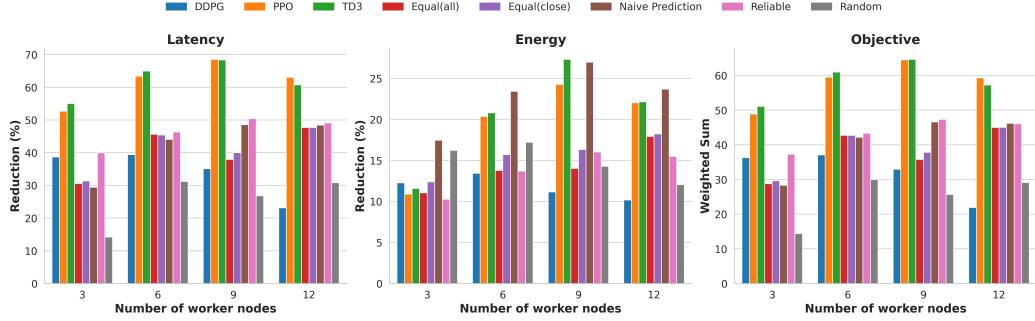


Figure 9: Performance Comparison ($200 \times 200m^2$)

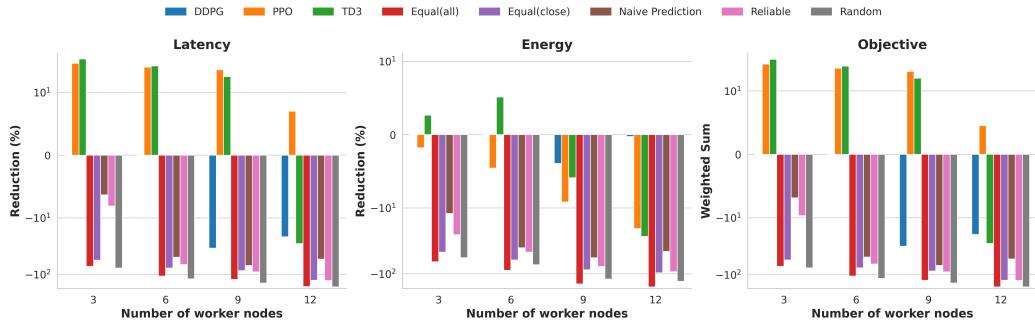


Figure 10: Performance Comparison ($300 \times 300m^2$)

389 master satisfying $d_i < 100m$; (3) **Naive Prediction** that selects offloadees
 390 based on predicted future trajectories and assigns sub-tasks of equal size to
 391 these offloadees. Specifically, it predicts each UAV's trajectory for the next
 392 20 steps, assuming the UAVs maintain their current velocity and ignoring
 393 potential collisions. It then calculates the percentage q_i of predicted UAV
 394 positions that remain within 200m of the master ($d_i < 200m$). UAVs with
 395 $q_i \geq 40\%$ are selected as offloadees; (4) **Reliable** that allocates tasks based
 396 on a reliability score defined as reliability = $q_i \psi_i f_i$. It picks the top $\lceil \frac{N+1}{2} \rceil$
 397 UAVs with the highest reliability scores and assigns tasks to these UAVs
 398 proportionally to their reliability scores; (5) **Random** that randomly sample
 399 actions from the action space.

400 The results in Fig. 11 and Fig. 12 demonstrate the promising perfor-
 401 mance of our method, which achieves the highest latency reduction across
 402 all scenarios. When the area is $300 \times 300m^2$, the **Naive Prediction** ranks

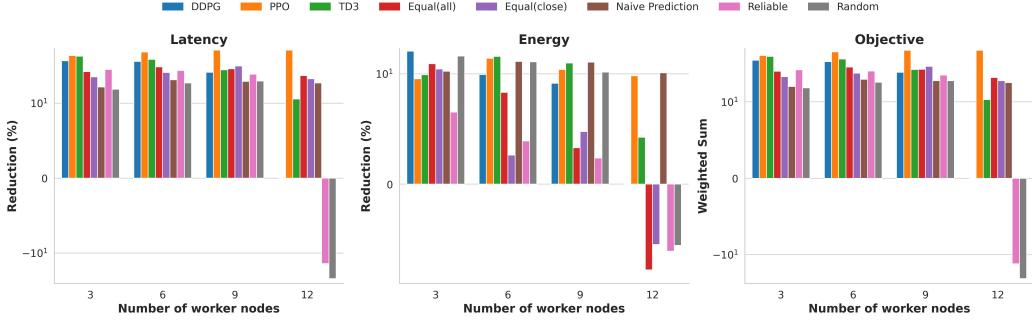


Figure 11: Performance Comparison ($200 \times 200m^2$)

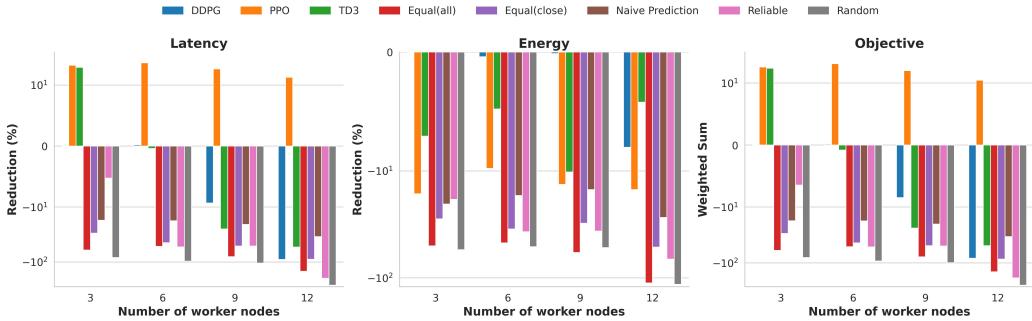


Figure 12: Performance Comparison ($300 \times 300m^2$)

403 second in scenarios with 3, 6, and 9 workers, while the **Equal (close)** ranks
 404 second in scenarios with 12 workers. When the area expands to $400 \times 400m^2$
 405 (Fig. 12), the **Equal (all)** and the **Random** provide no benefit in reducing
 406 latency; instead, they significantly delay task completion. Comparing Fig.
 407 11 and Fig. 12, we can observe that as the area increases for a fixed N , the
 408 performance of all methods degrades. This is due to the sparser airspace,
 409 which causes UAVs to be farther apart from each other, thereby increasing
 410 transmission latency and task completion time.

411 Appendix A. Example Appendix Section

412 Appendix text.

413 Example citation, See [4].

414 **References**

- 415 [1] H. Kurunathan, H. Huang, K. Li, W. Ni, E. Hossain, Machine learning-
416 aided operations and communications of unmanned aerial vehicles:
417 A contemporary survey, *IEEE Communications Surveys & Tutorials*
418 (2023).
- 419 [2] Y. Zeng, R. Zhang, T. J. Lim, Wireless communications with unmanned
420 aerial vehicles: Opportunities and challenges, *IEEE Communications*
421 magazine 54 (5) (2016) 36–42.
- 422 [3] B. Liu, W. Zhang, W. Chen, H. Huang, S. Guo, Online computation
423 offloading and traffic routing for uav swarms in edge-cloud computing,
424 *IEEE Transactions on Vehicular Technology* 69 (8) (2020) 8777–8791.
- 425 [4] Z. Bai, Y. Lin, Y. Cao, W. Wang, Delay-aware cooperative task offload-
426 ing for multi-uav enabled edge-cloud computing, *IEEE Transactions on*
427 *Mobile Computing* (2022).
- 428 [5] M. Satyanarayanan, The emergence of edge computing, *Computer* 50 (1)
429 (2017) 30–39.
- 430 [6] Q. Hu, Y. Cai, G. Yu, Z. Qin, M. Zhao, G. Y. Li, Joint offloading and
431 trajectory design for uav-enabled mobile edge computing systems, *IEEE*
432 *Internet of Things Journal* 6 (2) (2018) 1879–1892.
- 433 [7] Y. Miao, K. Hwang, D. Wu, Y. Hao, M. Chen, Drone swarm path plan-
434 ning for mobile edge computing in industrial internet of things, *IEEE*
435 *Transactions on Industrial Informatics* (2022).
- 436 [8] K. Lu, J. Xie, Y. Wan, S. Fu, Toward uav-based airborne computing,
437 *IEEE Wireless Communications* 26 (6) (2019) 172–179.
- 438 [9] H. Zhang, B. Wang, R. Wu, J. Xie, Y. Wan, S. Fu, K. Lu, Exploring net-
439 worked airborne computing: A comprehensive approach with advanced
440 simulator and hardware testbed, *Unmanned Systems* (2023).
- 441 [10] B. Wang, J. Xie, K. Lu, Y. Wan, S. Fu, Learning and batch-processing
442 based coded computation with mobility awareness for networked air-
443 borne computing, *IEEE Transactions on Vehicular Technology* (2022).

- 444 [11] E. Shtaiwi, A. Abdelhadi, H. Li, Z. Han, H. V. Poor, Orthogonal time
445 frequency space for integrated sensing and communication: A survey,
446 arXiv preprint arXiv:2402.09637 (2024).
- 447 [12] W. Lu, P. Si, Y. Gao, H. Han, Z. Liu, Y. Wu, Y. Gong, Trajectory and
448 resource optimization in ofdm-based uav-powered iot network, IEEE
449 Transactions on Green Communications and Networking 5 (3) (2021)
450 1259–1270.
- 451 [13] X. Guan, Y. Huang, Q. Shi, Joint subcarrier and power allocation for
452 multi-uav systems, China Communications 16 (1) (2019) 47–56.
- 453 [14] E. M. Royer, P. M. Melliar-Smith, L. E. Moser, An analysis of the
454 optimum node density for ad hoc mobile networks, in: ICC 2001. IEEE
455 International Conference on Communications. Conference Record (Cat.
456 No. 01CH37240), Vol. 3, IEEE, 2001, pp. 857–861.
- 457 [15] S. Fujimoto, H. van Hoof, D. Meger, Addressing Function Approxima-
458 tion Error in Actor-Critic Methods (Oct. 2018). arXiv:1802.09477.
- 459 [16] K. Cheng, Y. Teng, W. Sun, A. Liu, X. Wang, Energy-efficient joint
460 offloading and wireless resource allocation strategy in multi-mec server
461 systems, in: 2018 IEEE international conference on communications
462 (ICC), IEEE, 2018, pp. 1–6.
- 463 [17] F. Pervez, A. Sultana, C. Yang, L. Zhao, Energy and latency efficient
464 joint communication and computation optimization in a multi-uav as-
465 sisted mec network, IEEE Transactions on Wireless Communications
466 (2023).
- 467 [18] A. Goldsmith, Wireless communications, Cambridge university press,
468 2005.
- 469 [19] X. Zhang, J. Xie, Drl-based task offloading for networked uavs with
470 random mobility and collision avoidance, in: 2024 20th International
471 Conference on Wireless and Mobile Computing, Networking and Com-
472 munications (WiMob), IEEE, 2024, pp. 514–519.
- 473 [20] D. H. Choi, S. H. Kim, D. K. Sung, Energy-efficient maneuvering
474 and communication of a single uav-based relay, IEEE Transactions on
475 Aerospace and Electronic Systems 50 (3) (2014) 2320–2327.

- 476 [21] Y. Wan, K. Namuduri, Y. Zhou, D. He, S. Fu, A smooth-turn mobility
477 model for airborne networks, in: Proceedings of the first ACM MobiHoc
478 workshop on Airborne Networks and Communications, 2012, pp. 25–30.
- 479 [22] D. S. Lakew, U. Sa'ad, N.-N. Dao, W. Na, S. Cho, Routing in flying ad
480 hoc networks: A comprehensive survey, IEEE Communications Surveys
481 & Tutorials 22 (2) (2020) 1071–1120.
- 482 [23] N. T. Hoa, N. C. Luong, D. Van Le, D. Niyato, et al., Deep reinforcement
483 learning for multi-hop offloading in uav-assisted edge computing, IEEE
484 Transactions on Vehicular Technology (2023).
- 485 [24] C. M. Bishop, N. M. Nasrabadi, Pattern recognition and machine learn-
486 ing, Vol. 4, Springer, 2006.
- 487 [25] D. E. Rumelhart, G. E. Hinton, R. J. Williams, Learning representations
488 by back-propagating errors, nature 323 (6088) (1986) 533–536.
- 489 [26] K. He, X. Zhang, S. Ren, J. Sun, Delving deep into rectifiers: Surpassing
490 human-level performance on imagenet classification, in: Proceedings of
491 the IEEE international conference on computer vision, 2015, pp. 1026–
492 1034.