

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,c,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*

^c*Department of Civil and Environmental Engineering, University of Pittsburgh, 3700 O'Hara Street, Pittsburgh, 15261, Pennsylvania, 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

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

*Corresponding author

Email addresses: xiz166@ucsd.edu (Xixin Zhang), doj14@pitt.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.

approach significantly speeds up task execution compared to existing methods.

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

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

In this paper, we investigate a more common yet challenging NAC scenario where all UAVs, including both offloaders and offloaders, move randomly during task execution while actively avoiding collisions. None of the UAVs have prior knowledge of the environment or system models, and their movement patterns or future trajectories are not shared among each other. Additionally, we generalize computation tasks as any functions or operations that can be partitioned into arbitrary subtasks for parallel computation. To model UAV movement, we extend the traditional Random Direction model [11], originally designed for individual entities, to capture collision avoidance interactions among multiple UAVs. Furthermore, we formulate a nonlinear optimization to optimize task allocation and develop a DRL algorithm based on the Twin Delayed Deep Deterministic Policy Gradient (TD3)[12] to solve it. We evaluate the performance of the proposed method through extensive comparative simulation studies, which demonstrate its promising

70 performance.

71 In the rest of this paper, Sec. 2 details the system models and formulates
72 the optimization problem. Sec. 3 describes the proposed DRL algorithm.
73 In Sec. 4, simulation results are presented and discussed. We conclude in
74 Sec. ??.

75 2. System Models and Problem Formulation

76 In this section, we first introduce the system models used to construct the
77 simulated environment, which are unknown to the UAVs. Then we formulate
78 the problem to be solved.

79 2.1. System Models

80 Consider a group of $N + 1$ heterogeneous UAVs with varying physical
81 configurations, indexed as $i \in \mathcal{N} = \{0, 1, 2, \dots, N\}$. Each UAV is equipped
82 with computing and communication modules, enabling resource sharing and
83 onboard computation. Their computing, communication, and mobility char-
84 acteristics can be modeled as follows.

85 2.1.1. Computing Model

86 We describe the computing capability of each UAV i as CPU cycle fre-
87 quency f_i in Hz. For a general computing task k , its input data size is S_k
88 (bits), and its required computation intensity is ξ_k (cycles/bit)[13]. The total
89 CPU cycles required to compute task k is hence $\xi_k S_k$ and the time required
90 for UAV i to execute this task is

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

91 2.1.2. Communication Model

92 Denote the distance between UAV i and UAV j as d_{ij} . The UAV-to-UAV
93 links are typically Line of Sight (LoS), with propagation speed approaching
94 the speed of light. Hence, the transmission latency can be approximated
95 using the transmission time. Here, we model the transmission rate (bits/s)
96 based on the Simplified Path Loss Model [14] as follows

$$\nu_{ij} = B \log_2 \left(1 + \frac{G(d_r/d_{ij})^\theta \psi_i}{N_0 B} \right) \quad (2)$$

where B is the bandwidth (Hz) of the channel, d_r is the reference distance (meter), G is the unitless constant equal to the path gain of the distance d_r , θ is the path loss exponent, ψ_i is the transmitted power (mW) and N_0 is the noise power spectral density (dBm/Hz). The overall transmission time is as follows

$$T_{k,ij}^{trans} = \frac{S_k}{\nu_{ij}} \quad (3)$$

2.1.3. Mobility Model

We assume the UAVs fly at the same altitude. Therefore, the position 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 constraints $0 \leq x_i(t) \leq W$, $0 \leq y_i(t) \leq W$, such that the position is bounded within an area of $W \times W (m^2)$. To model its movement, we adopt the Random Direction (RD) model [11], which has been widely used for describing UAVs, particularly multirotor drones [15]. Given initial position $\mathbf{p}_i(0) = (x_i(0), y_i(0))$, UAV i randomly picks a constant velocity, where the magnitude ranges uniformly between 0 and $v_{max} \in \mathbb{R}$, and the direction is uniformly distributed across 2π . The UAV then moves in a straight line for a duration randomly selected from an exponential distribution with parameter λ . Once this duration is completed, the UAV chooses another velocity and duration, repeating the process. Suppose at the m -th instance when UAV i changes its velocity, it selects a new velocity $\mathbf{v}_{i,m} \in \mathbb{R}^2$ and duration $\delta_{i,m} \in \mathbb{R}$. We define the start time of the m -th instance as $t_m = \sum_{l=-1}^{m-1} \delta_{i,l}$, with $\delta_{i,-1} = 0$. The UAV i 's position at time t , where $t_m \leq t < t_{m+1}$, can be represented as follows

$$\mathbf{p}_i(t) = \mathbf{p}_i(0) + \sum_{l=-1}^{m-1} \delta_{i,l} \mathbf{v}_{i,l} + (t - t_m) \mathbf{v}_{i,m} \quad (4)$$

The traditional RD model is designed to describe the mobility of a single entity. However, in multi-UAV systems, the mobility of UAVs can change to avoid collisions. This necessitates the incorporation of collision avoidance mechanisms into the RD model. Here, we assume that each UAV i will turn around and move in the opposite direction without changing speed until the current duration is completed when its distance to any other UAV j falls below a threshold D_i . Note that if both UAVs 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

128 area. The mobility of each UAV i with collision avoidance can then be
 129 described as

$$\begin{aligned} \mathbf{p}_i(t) = \mathbf{p}_i(0) + \sum_{l=-1}^{m-1} (\delta_{i,l}^o - \delta_{i,l}^c) \mathbf{v}_{i,l} \\ + (\tilde{\delta}_{i,m}^o - \tilde{\delta}_{i,m}^c) \mathbf{v}_{i,m} \end{aligned} \quad (5)$$

130 where $0 \leq \delta_{i,l}^o \leq \delta_{i,l}$ is the time spent moving at velocity $\mathbf{v}_{i,l}$ in the l -th
 131 instance before triggering collision avoidance and $\delta_{i,l}^c = \delta_{i,l} - \delta_{i,l}^o$. Likewise,
 132 $0 \leq \tilde{\delta}_{i,m}^o \leq t - t_m$ is the time spent moving at velocity $\mathbf{v}_{i,m}$ before collision
 133 avoidance in the m -th instance, and $\tilde{\delta}_m^c = (t - t_m) - \tilde{\delta}_m^o$.

134 2.2. Problem Formulation

135 Without loss of generality, we let UAV $i = 0$ be the master (or offloader)
 136 and treat the remaining N UAVs as potential offloaders with idle computing
 137 resources. Suppose a sequence of computing tasks $\mathcal{K} = \{1, 2, \dots, K\}$ is gener-
 138 ated at the master, and each task k can be divided into $L_k \in \mathbb{Z}^+$ atomic tasks
 139 of size $\ell_k = \frac{S_k}{L_k}$, which can be computed in parallel. To minimize the total
 140 task completion time, the master aims to optimally allocate the L_k atomic
 141 tasks among all the UAVs, including itself. Suppose the workload offloaded
 142 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
 143 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
 144 workload offloaded to UAV i , then $c_i^k = 0$.

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

$$T_k = \max_i T_{k,i} \quad (6)$$

148 For simplicity, we assume the result size is small and the latency of sending
 149 it back is negligible, as often assumed in existing works [16]. Therefore, $T_{k,i}$
 150 can be expressed as

$$T_{k,i} = T_{k,i}^{comp} + T_{k,i}^{trans} \quad (7)$$

151 According to the computing and communication models described in the
 152 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

$$\begin{cases} \int_0^{T_{k,i}^{trans}} \nu_{0i}(t) dt = c_i^k \ell_k & \text{if } i \neq 0 \\ T_{k,i}^{trans} = 0 & \text{if } i = 0 \end{cases} \quad (8)$$

153 The problem can then be mathematically formulated as follows

$$\underset{c_i^k, \forall i \in \mathcal{N}, k \in \mathcal{K}}{\text{Minimize}} \quad \sum_{k=1}^K T_k \quad (9a)$$

$$\text{subject to} \quad \sum_{i=0}^N c_i^k \ell_k = S_k \quad \forall k \in \mathcal{K} \quad (9b)$$

$$c_i^k \in \mathbb{Z}, 0 \leq c_i^k \leq L_k, \forall i \in \mathcal{N}, k \in \mathcal{K} \quad (9c)$$

154 3. Deep Reinforcement Learning Solution

155 To solve the optimization problem formulated in the previous section, the
 156 greatest challenge lies in the unknown relationship between T_k and the deci-
 157 sion variables c_i^k , as UAVs lack knowledge of the system models. Additionally,
 158 the randomness of UAVs' mobility presents another significant challenge. In
 159 this section, we introduce our model-free DRL-based algorithm, a variant of
 160 the Twin Delayed Deep Deterministic policy gradient algorithm (TD3) [12],
 161 to address these challenges.

162 3.1. Markov Decision Process

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

168 3.1.1. State

169 We assume the master agent only has access to all UAVs' positions
 170 $\mathbf{p}_i(t), \forall i \in \mathcal{N}$ at each time t , and the prior knowledge about their config-
 171 urations, *i.e.* idle computing resources f_i , and communication capabilities
 172 characterized by transmitted powers ψ_i . We discretize the time into steps of
 173 length Δt . Let t^k denote the start time of task k , the state at t^k is defined as
 174 the combination of the task size S_k , all UAVs' historic trajectories and con-
 175 figurations $s_k = [S_k, (\boldsymbol{\rho}_0(k), f_0, \psi_0), (\boldsymbol{\rho}_1(k), f_1, \psi_1), \dots, (\boldsymbol{\rho}_N(k), f_N, \psi_N)] \in \mathcal{S}$,
 176 where $\boldsymbol{\rho}_i(k) = [\mathbf{p}_i(t^k - M\Delta t), \mathbf{p}_i(t^k - (M-1)\Delta t), \dots, \mathbf{p}_i(t^k)] \in \mathbb{R}^{2 \times (M+1)}$ is
 177 UAV i 's trajectory consisting of its most recent M -step positions (including
 178 the start point).

179 3.1.2. Action

180 Every time a task k arrives, the master agent partitions it into sub-tasks
 181 of varying amounts of atomic tasks and offloads them to different UAVs.
 182 Let $\mu_i^k \geq 0$ denote the portion of task k offloaded to UAV $i \in \mathcal{N}$, such
 183 that $\sum_{i=0}^N \mu_i^k = 1$. The action taken by the master at time t^k is defined as
 184 $a_k = [\mu_0^k, \mu_1^k, \dots, \mu_N^k] \in \mathcal{A}$, where \mathcal{A} is the space of N -simplex. The workload
 185 offloaded to UAV i is then computed by $c_i^k = \text{round}(\frac{\mu_i^k S_k}{\ell_k})$.

186 3.1.3. Reward

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

$$r(s_k, a_k) = \frac{S_k}{\sum_{k \in \mathcal{K}} S_k} \cdot \frac{\tilde{T}_k - T_k}{\tilde{T}_k} \quad (10)$$

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

192 3.1.4. Transition

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

$$s_{k+1} \sim P(s_{k+1} | s_k, a_k; v_{max}, \lambda) \quad (11)$$

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

199 3.2. TD3-based Offloading

200 TD3 [12] is an advanced actor-critic algorithm designed for continuous
 201 action spaces in DRL. Here, we introduce a variant of the TD3 algorithm to
 202 enable the master UAV to identify trustworthy offloaders and optimize task
 203 allocation.

204 3.2.1. Actor

205 The behavior of the master agent is defined by a policy function $\pi : \mathcal{S} \rightarrow$
 206 \mathcal{A} , which maps each state $s \in \mathcal{S}$ to a continuous action $a \in \mathcal{A}$. The goal of the
 207 agent is to determine the optimal policy π^* , which maximizes the expected
 208 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
 209 accumulative discounted reward and γ is the discount factor.

210 To approximate the policy function, TD3 [12] utilizes a neural network
 211 parameterized by ϕ , denoted as π_ϕ , which directly maps the state space $s \in \mathcal{S}$
 212 to the action space \mathcal{A} . To satisfy the constraints in (9b), we adjust the actor
 213 π_ϕ as follows. By including all UAVs' historic trajectories in the state, the
 214 neural network can learn the movement patterns of each UAV and their inter-
 215 actions. Moreover, the network adjusts weights to prioritize UAVs that have
 216 more resources to share and are more likely to remain close to the master.
 217 Therefore, the neural network's output logits $\mathbf{z}^k \in \mathbb{R}^{N+1}$ can be interpreted
 218 as the reliability score of each UAV. We then use the Softmax function to
 219 decide the offloading portions according to UAVs' reliability scores as follows

$$\mu_i^k = \text{Softmax}(z_i^k) = \frac{e^{z_i^k}}{\sum_{j=0}^N e^{z_j^k}} \quad (12)$$

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

221 To estimate the parameters ϕ , we apply off-policy learning to enhance
 222 training stability. Moreover, to strengthen the robustness of the learned
 223 policy function against variance and prevent overfitting to narrow peaks of
 224 action values, we add random noise to the actions of the target actor $\pi_{\phi'}$
 225 parameterized by ϕ' as follows [12]

$$\begin{aligned} \tilde{\mu}_i^k &= (1 - \beta)\mu_i^k + \beta\epsilon_i \\ \epsilon &\sim \text{Dir}(\alpha_{policy}) \end{aligned} \quad (13)$$

226 where β is the weight of noise ϵ_i , and $\epsilon_i \in \mathbb{R}^{N+1}$ is sampled from the Dirichlet
 227 distribution [17] with a concentration parameter α_{policy} that is uniform across
 228 all elements. Of note, the Dirichlet distribution guarantees that the actions
 229 remain within the space of N -simplex.

230 3.2.2. Critic

231 Given the state s_k and the action a_k taken, the state-action value func-
 232 tion Q^π provides the expected return when following policy π thereafter, i.e.,

233 $Q^\pi(s_k, a_k) = \mathbb{E}_{s_i >_k \sim P, a_i >_k \sim \pi} [R_k | s_k, a_k]$. To approximate the critic Q^π , neural
 234 networks are also used. Following TD3[12], we define two primary critic net-
 235 works Q_{θ_1} and Q_{θ_2} and two target critic networks $Q_{\theta'_1}$ and $Q_{\theta'_2}$ parameterized
 236 by $\theta_1, \theta_2, \theta'_1$ and θ'_2 , respectively.

237 3.2.3. Training

238 As shown in Alg. 1, the training starts by initializing all the parameters
 239 and the replay buffer \mathcal{D} , which stores transition samples (Lines 1-3). At
 240 each training iteration u , the master agent interacts with the environment
 241 to collect transition data and store them in the buffer \mathcal{D} until the number
 242 of collected transitions exceeds H (Lines 5-6). After that, the agent utilizes
 243 a batch \mathcal{B} sampled from the replay buffer to update the parameters of the
 244 critics (Lines 9-11) and actor (Lines 12-13). Particularly, the parameters θ_i ,
 245 $i \in \{1, 2\}$, of each primary critic is updated by minimizing the following loss
 246 over the sampled batch

$$l_i = \frac{1}{|\mathcal{B}|} \sum_{(s,a,r,s') \in \mathcal{B}} (y - Q_{\theta_i}(s, a))^2 \quad (14)$$

247 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-
 248 larized action defined in (13).

249 A lower update frequency is necessary for the policy function compared to
 250 the value function, otherwise, it may lead to divergence. Hence, we update
 251 the policy function every d iterations by maximizing the expected return
 252 in the direction of the batch gradient of the policy, where the gradient is
 253 computed by

$$\nabla_\phi J = \frac{1}{|\mathcal{B}|} \sum_{(s,a,r,s') \in \mathcal{B}} [\nabla_a Q_{\theta_1}(s, a)|_{a=\pi_\phi(s)} \nabla_\phi \pi_\phi(s)] \quad (15)$$

254 Also, the parameters of the target networks are gradually adjusted towards
 255 the weights of the primary networks through weighted soft updates every d
 256 iteration as follows

$$\begin{aligned} \theta'_i &\leftarrow \theta_i + (1 - \tau)(\theta'_i - \theta_i), i = 1, 2 \\ \phi' &\leftarrow \phi + (1 - \tau)(\phi' - \phi) \end{aligned} \quad (16)$$

257 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 $\epsilon \sim \text{Dir}(\alpha_{\text{explore}})$, 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: $\tilde{a}' \leftarrow (1 - \beta)\pi_{\phi'}(s') + \beta\epsilon$, $\epsilon \sim \text{Dir}(\alpha_{\text{policy}})$
 - 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**
-

258 4. Experiments

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

261 4.1. Environment Setting

262 We evaluate our method in simulated scenarios with one master UAV
263 and $N = 3, 6, 9$, or 12 offloadee UAVs respectively. For the multi-UAV RD
264 mobility model, we set its parameters as $v_{max} = 20$ m/s, $\lambda = \frac{1}{15}$, and $D_i = 5$,
265 $\forall i \in \mathcal{N}$. The length of each discrete time step is set to $\Delta t = 1$ second. We
266 also vary the size of the flying zone and consider two sizes, $W = 300$ m and
267 $W = 400$ m. The computing power f_i of each UAV i is randomly configured
268 by selecting values from the range of $[1, 1.6]$ Ghz. The task list consists of
269 $K = 25$ tasks that can be locally completed in $\tilde{T}_k \in [20, 60]$ seconds by the
270 master UAV. All tasks can be divided into $L_k = 1000$ atomic tasks. The
271 computation intensity is set to be the same $\xi_k = 10^6$ cycles/kB for all tasks,
272 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
273 UAV swarm utilizing 5G Wi-Fi communication for data transmission, we set
274 the bandwidth $B = 40$ MHz, relative distance $d_r = 1$ meter, path gain
275 $G = 40$, path loss exponent $\theta = 4$, and noise power spectral density $N_0 =$
276 -174 dBm/Hz. The transmitted power ψ_i of the master UAV to each UAV
277 $i \in \mathcal{N} \setminus \{0\}$ varies from 80 mW to 120 mW.

278 4.2. Training Performance

279 We employ a 3-layer Multilayer Perceptrons (MLP)[18] architecture for
280 both the actor and critic networks. The actor network takes observations as
281 input, whose dimension is based on the number of computing nodes $N + 1$,
282 and outputs actions of dimension N . Each UAV trajectory has a length
283 of $M + 1$, where $M = 20$. The critic network takes the concatenation of
284 observations and actions as input and outputs a scalar representing the state
285 action value. The width of the hidden layer in both networks is set to twice
286 the dimension of the input. All parameters are initialized using Kaiming
287 initialization[19].

288 The learning rates for the actor and critic networks are set to 0.0001 and
289 0.0002, respectively. The threshold H is set to 1000 and the batch size is
290 $|\mathcal{B}| = 512$. The gradient norm is clipped between 0 and 0.2. The exploration
291 and policy noise are sampled from a Dirichlet distribution with parameters
292 $\alpha_{explore} = 0.1$ and $\alpha_{policy} = 0.99$, respectively. The noise weight is $\beta = 0.1$

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

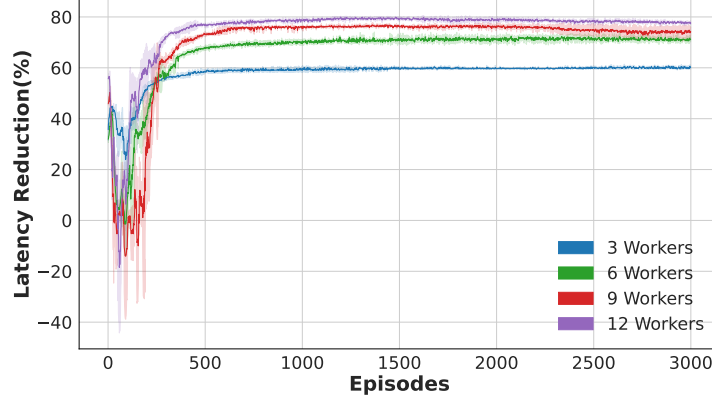


Figure 1: Learning Curve ($300 \times 300m^2$)

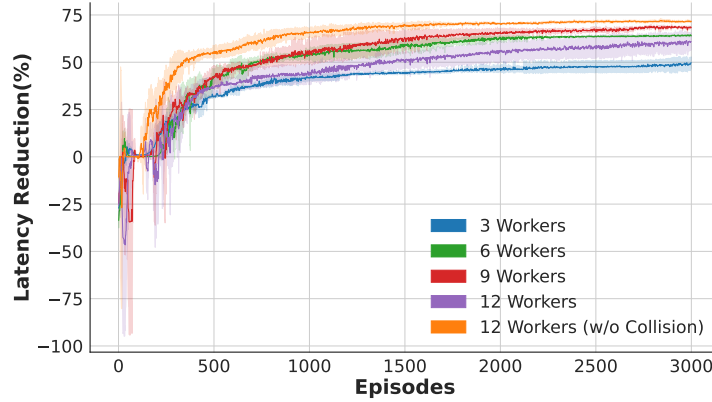


Figure 2: Learning Curve ($400 \times 400m^2$)

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. 1 and Fig. 2. 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.,

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

295 The two figures demonstrate that our method converges after training and
 296 shows promising stability. When UAVs are restricted to an area of $300 \times$
 297 $300m^2$ (Fig. 1), increasing the number of potential offloaders (workers) re-
 298 duces task completion time. The same phenomenon is observed when the
 299 flying zone expands to $400 \times 400m^2$, except when the number of potential
 300 workers increases to $N = 12$. The performance degradation at $N = 12$ may
 301 be due to increased collision avoidance maneuvers, which make UAVs' mobil-
 302 ity more uncertain and harder to learn and predict. Intuitively, the master
 303 agent tends to share workload with UAVs that are more likely to remain
 304 nearby throughout task execution, as indicated by a higher reliability score
 305 z_i^k . Therefore, when UAV mobility becomes more uncertain, fewer workers
 306 are selected as offloaders due to reduced reliability. This is confirmed by the
 307 results shown in Fig. 2, where performance improves when collision avoid-
 308 ance mechanisms are not in place. It also indicates that the task completion
 309 time cannot be infinitely reduced by continuously adding more UAVs to the
 310 region.

311 4.3. Comparison Studies

312 To evaluate the performance of the proposed method, we compare it
 313 with five benchmarks, including (1) **Equal (all)** that equally divides and
 314 distributes tasks to all UAVs; (2) **Equal (close)** that equally partitions
 315 and distributes tasks to UAVs that are close enough with distance to the
 316 master satisfying $d_i < 100m$; (3) **Naive Prediction** that selects offloaders
 317 based on predicted future trajectories and assigns sub-tasks of equal size to
 318 these offloaders. Specifically, it predicts each UAV's trajectory for the next
 319 20 steps, assuming the UAVs maintain their current velocity and ignoring
 320 potential collisions. It then calculates the percentage q_i of predicted UAV
 321 positions that remain within 200m of the master ($d_i < 200m$). UAVs with
 322 $q_i \geq 40\%$ are selected as offloaders; (4) **Reliable** that allocates tasks based
 323 on a reliability score defined as $\text{reliability} = q_i \psi_i f_i$. It picks the top $\lceil \frac{N+1}{2} \rceil$
 324 UAVs with the highest reliability scores and assigns tasks to these UAVs
 325 proportionally to their reliability scores; (5) **Random** that randomly samples
 326 actions from the action space.

327 The results in Fig. 3 and Fig. 4 demonstrate the promising performance
 328 of our method, which achieves the highest latency reduction across all sce-
 329 narios. When the area is $300 \times 300m^2$, the **Naive Prediction** ranks second
 330 in scenarios with 3, 6, and 9 workers, while the **Equal (close)** ranks second
 331 in scenarios with 12 workers. When the area expands to $400 \times 400m^2$ (Fig.

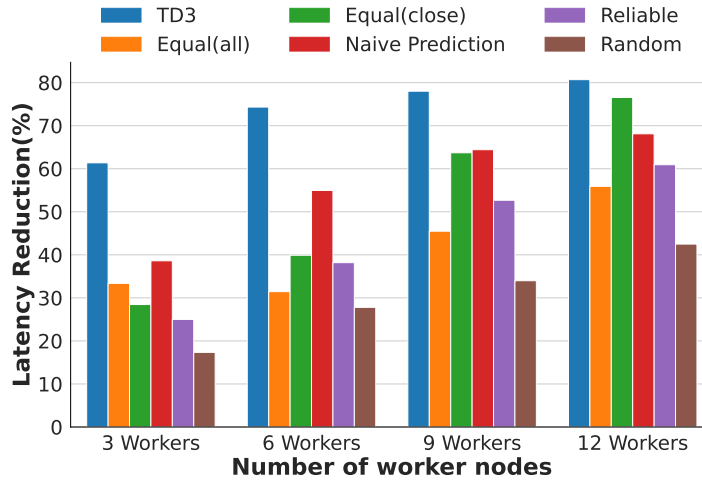


Figure 3: Performance Comparison ($300 \times 300m^2$)

4), the **Equal (all)** and the **Random** provide no benefit in reducing latency; instead, they significantly delay task completion. Comparing Fig. 3 and Fig. 4, 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 transmission latency and task completion time.

Appendix A. Example Appendix Section

Appendix text.
Example citation, See [4].

References

- [1] H. Kurunathan, H. Huang, K. Li, W. Ni, E. Hossain, Machine learning-aided operations and communications of unmanned aerial vehicles: A contemporary survey, *IEEE Communications Surveys & Tutorials* (2023).
- [2] Y. Zeng, R. Zhang, T. J. Lim, Wireless communications with unmanned aerial vehicles: Opportunities and challenges, *IEEE Communications magazine* 54 (5) (2016) 36–42.

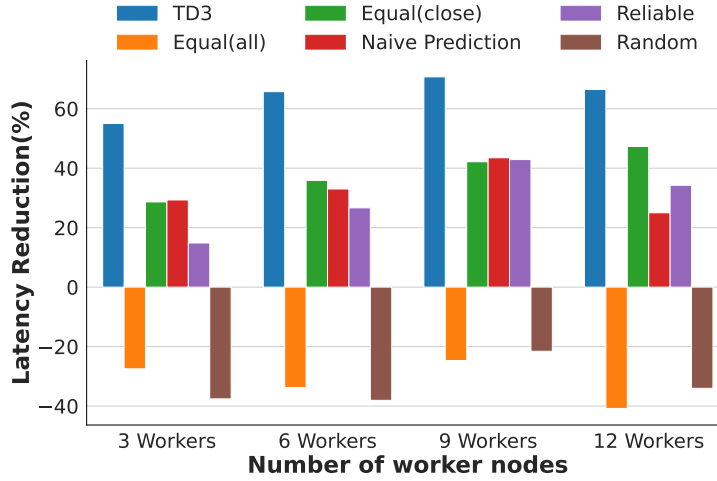


Figure 4: Performance Comparison ($400 \times 400m^2$)

- 349 [3] B. Liu, W. Zhang, W. Chen, H. Huang, S. Guo, Online computation
 350 offloading and traffic routing for uav swarms in edge-cloud computing,
 351 IEEE Transactions on Vehicular Technology 69 (8) (2020) 8777–8791.
- 352 [4] Z. Bai, Y. Lin, Y. Cao, W. Wang, Delay-aware cooperative task offload-
 353 ing for multi-uav enabled edge-cloud computing, IEEE Transactions on
 354 Mobile Computing (2022).
- 355 [5] M. Satyanarayanan, The emergence of edge computing, Computer 50 (1)
 356 (2017) 30–39.
- 357 [6] Q. Hu, Y. Cai, G. Yu, Z. Qin, M. Zhao, G. Y. Li, Joint offloading and
 358 trajectory design for uav-enabled mobile edge computing systems, IEEE
 359 Internet of Things Journal 6 (2) (2018) 1879–1892.
- 360 [7] Y. Miao, K. Hwang, D. Wu, Y. Hao, M. Chen, Drone swarm path plan-
 361 ning for mobile edge computing in industrial internet of things, IEEE
 362 Transactions on Industrial Informatics (2022).
- 363 [8] K. Lu, J. Xie, Y. Wan, S. Fu, Toward uav-based airborne computing,
 364 IEEE Wireless Communications 26 (6) (2019) 172–179.
- 365 [9] H. Zhang, B. Wang, R. Wu, J. Xie, Y. Wan, S. Fu, K. Lu, Exploring net-
 366 worked airborne computing: A comprehensive approach with advanced
 367 simulator and hardware testbed, Unmanned Systems (2023).

- 368 [10] B. Wang, J. Xie, K. Lu, Y. Wan, S. Fu, Learning and batch-processing
369 based coded computation with mobility awareness for networked air-
370 borne computing, *IEEE Transactions on Vehicular Technology* (2022).
- 371 [11] E. M. Royer, P. M. Melliar-Smith, L. E. Moser, An analysis of the
372 optimum node density for ad hoc mobile networks, in: *ICC 2001. IEEE*
373 *International Conference on Communications. Conference Record (Cat.*
374 *No. 01CH37240)*, Vol. 3, IEEE, 2001, pp. 857–861.
- 375 [12] S. Fujimoto, H. van Hoof, D. Meger, Addressing Function Approxima-
376 tion Error in Actor-Critic Methods (Oct. 2018). [arXiv:1802.09477](https://arxiv.org/abs/1802.09477).
- 377 [13] K. Cheng, Y. Teng, W. Sun, A. Liu, X. Wang, Energy-efficient joint
378 offloading and wireless resource allocation strategy in multi-mec server
379 systems, in: *2018 IEEE international conference on communications*
380 *(ICC)*, IEEE, 2018, pp. 1–6.
- 381 [14] A. Goldsmith, *Wireless communications*, Cambridge university press,
382 2005.
- 383 [15] D. S. Lakew, U. Sa’ad, N.-N. Dao, W. Na, S. Cho, Routing in flying ad
384 hoc networks: A comprehensive survey, *IEEE Communications Surveys*
385 *& Tutorials* 22 (2) (2020) 1071–1120.
- 386 [16] N. T. Hoa, N. C. Luong, D. Van Le, D. Niyato, et al., Deep reinforcement
387 learning for multi-hop offloading in uav-assisted edge computing, *IEEE*
388 *Transactions on Vehicular Technology* (2023).
- 389 [17] C. M. Bishop, N. M. Nasrabadi, *Pattern recognition and machine learn-*
390 *ing*, Vol. 4, Springer, 2006.
- 391 [18] D. E. Rumelhart, G. E. Hinton, R. J. Williams, Learning representations
392 by back-propagating errors, *nature* 323 (6088) (1986) 533–536.
- 393 [19] K. He, X. Zhang, S. Ren, J. Sun, Delving deep into rectifiers: Surpassing
394 human-level performance on imagenet classification, in: *Proceedings of*
395 *the IEEE international conference on computer vision*, 2015, pp. 1026–
396 1034.