# Securing the Skies: An IRS-Assisted AoI-Aware Secure Multi-UAV System with Efficient Task Offloading

[1] Joshi Poorvi, [2]Alakesh Kalita, [3]Mohan Gurusamy

[1,3]Electrical and Computer Engineering, National University of Singapore, Singapore

[2]ISTD Pillar, Singapore University of Technology and Design, Singapore

[1]e1144005@u.nus.edu, [2]alakesh_kalita@sutd.edu.sg, [3]gmohan@nus.edu.sg

*Abstract*—**Unmanned Aerial Vehicles (UAVs) are integral in various sectors like agriculture, surveillance, and logistics, driven by advancements in 5G. However, existing research lacks a comprehensive approach addressing both data freshness and security concerns. In this paper, we address the intricate challenges of data freshness, and security, especially in the context of eavesdropping and jamming in modern UAV networks. Our framework incorporates exponential AoI metrics and emphasizes secrecy rate to tackle eavesdropping and jamming threats. We introduce a transformer-enhanced Deep Reinforcement Learning (DRL) approach to optimize task offloading processes. Comparative analysis with existing algorithms showcases the superiority of our scheme, indicating its promising advancements in UAV network management.**

*Index Terms*—**Unmanned Aerial Vehicles, Age of Information, Intelligent Reflecting Surfaces, Deep Reinforcement Learning, Physical Layer Security, Data Freshness, Task Offloading.**

## I. Introduction

Unmanned Aerial Vehicles (UAVs), have seen increased utilization in agriculture, surveillance, logistics, and emergency response, facilitated by advancements in 5G networks, providing high-speed, low-latency communication. In UAV networks, maintaining data freshness is vital, especially in emergencies such as natural disasters or search and rescue operations. Data freshness refers to the timeliness of information collected by UAVs, closely linked to the Age of Information (AoI), measuring the duration between data acquisition and availability for decision-making. Recent studies underscore the importance of AoI in UAV networks. [1] emphasizes the importance of minimizing AoI for timely information collection and processing in UAV-aided Mobile Edge Computing Networks. Similarly, [2] investigates trajectory planning for multiple UAVs in IoT networks to minimize average AoI. Additionally, [3] utilizes AoI as a metric to assess temporal correlation in IoT data packets and proposes an AoI-energy-aware data collection scheme for UAV-assisted IoT networks. These studies underscore the significance of reducing overall AoI to maintain data freshness below a specified threshold. Furthermore, Security threats such as eavesdropping and jamming pose significant challenges to UAV network reliability and integrity. Various Physical Layer Security (PLS) methods, including multi-antenna relaying and artificial noise, aim to mitigate these threats [4]. These strategies aim to mitigate eavesdropping channels and reduce the interception of information by unauthorized receivers. Recent research has explored the use of Intelligent Reflecting Surfaces (IRS) to enhance PLS in UAV networks. IRS dynamically adjusts reflecting elements to improve signal reception by legitimate users while reducing unauthorized channel quality [5]. While previous studies have focused on individual optimization of UAV-IRS integration for AoI metrics [6], but simultaneous optimization of secrecy rate and AoI metrics in a unified framework remains unexplored. There's a need for a comprehensive framework addressing security challenges related to jamming, eavesdropping, and data freshness.

Our work focuses on optimizing an IRS-assisted AoI aware bi-layer multi-UAV system. A cooperative data sensing and transmission framework is developed where UAVs utilize radar signals to gather status information from users and securely relay it to the base station in the presence of eavesdroppers and jammers. The optimization challenge involves jointly considering the trajectory of the UAVs and the IRS reflection vector for PLS. The main contribution includes devising strategies to optimize the task offloading process within the proposed system, including designing non-overlapping trajectories for Computational-UAVs (C-UAVs) at the lower layer and IRS-aided UAVs (I-UAVs) operating at higher altitudes, focusing on collision avoidance and determining optimal beamforming vectors for I-UAVs. Further, a framework is introduced that incorporates exponential penalty-based AoI metrics and overall secrecy rate, aiming to enhance data freshness and security within network. At last, we used Multi-agent Deep Reinforcement Learning (DRL) supported by Gated Transformer architecture for joint optimization, facilitating efficient temporal modeling to optimize across multiple agents.

## II. System and Channel Model

In this study, we propose an IRS-assisted AoI aware bi-layer multi-UAV system designed to enhance the reliability of wireless networks catering to multiple User Equipment (UEs) in the timely execution of computation-intensive and delay-sensitive tasks. The proposed 3D system at time stamp
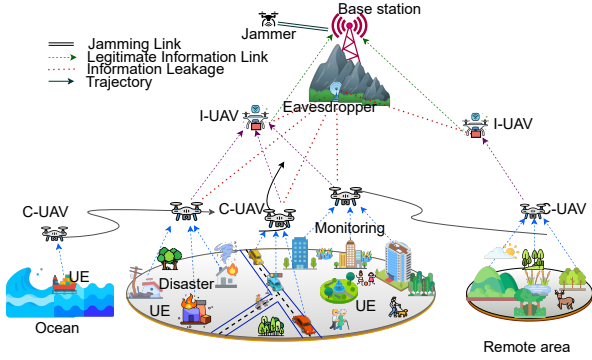
Fig. 1: The system model of proposed UAV network

## A. Transmission Protocol

The channel coefficients for various links are represented as follows: the $m^{th}$ user to $n^{th}$ C-UAV link ($h_{mn}^{UC}$), the $n^{th}$ C-UAV- eavesdropper link ($h_{n0}^{CE}$), the jammer-BS link ($h^{JB}$), the $n$th C-UAV - $p^{th}$ I-UAV link ($h_{np}^{CI} \in \mathbb{C}^{L \times 1}$), the jammer- $p$th I-UAV link ($h_{0p}^{JI} \in \mathbb{C}^{L \times 1}$), the $p$th I-UAV-BS link ($h_{p0}^{IB} \in \mathbb{C}^{1 \times L}$), and the $p$th I-UAV-E link ($h_{p0}^{IE} \in \mathbb{C}^{1 \times L}$). To model the channel properties we used Nakagami-m distribution [7]. Assuming accurate Channel State Information (CSI) for legitimate channels due to slow-varying nature. However, the lack of cooperation between the C-UAV and third-party nodes makes obtaining the CSI for illegitimate channels challenging, we assume partial knowledge of the CSI. So, CSI uncertainty for those channels are characterized using bounded CSI model [8]:

$$h^i = \hat{h}^i + \Delta h^i, \qquad (1)$$
$$\|\Delta h_i\| \leq \xi_{h,i}, \quad i \in \{CE, JB, JI, IE\}$$

where $\hat{h}_i$ represents the estimated CSI known at the C-UAV, and $\Delta h^i$ and is the unknown CSI error. Additionally, $\xi_{h,i}$ represents the levels of CSI uncertainty.

- **User to C-UAV transmission:** The received signal at $n$th C-UAV from $m$th UE can be expressed as

$$y_{C-UAV}^n = \sum_{m=1}^{M} \frac{h_{mn}^{UC}}{\sqrt{L_{mn}^{UC}}} \cdot x_{UE}^m + n_{C-UAV}^n \qquad (2)$$

where $x_{UE}^m$ is the transmitted symbol with energy $P_m$ from $m^{th}$ UE, and $n_{C-UAV}^n \sim N(0, \sigma_{C_n}^2)$ denotes the AWGN term with zero mean and variance $\sigma_{C_n}^2$. $L_{mn}^{UC} = A\|s_m^{UE}(t) - s_n^{C-UAV}(t)\|^{\alpha_{pl}}$ represents the path loss where $A$ is the constant associated with the signal frequency and transmission environment, and $\alpha_{pl}$ denotes the path loss exponent. The instantaneous received SNR at $n^{th}$ C-UAV receives signal from $m^{th}$ UE will be,

$$\Gamma_{nm}^{C-UAV}(t) = \frac{h_{mn}^{UC^2} \cdot P_m}{(\sigma_{C_n})^2 + \sum_{i=1, i \neq m}^{M} L_{in}^{UC} h_{in}^{UC^2} \cdot P_i} \qquad (3)$$

In the process of offloading tasks, it is assumed that the up link bandwidth $B_u$ is distributed equally among all UEs. As a result, the instantaneous data rate between $m^{th}$ UE, and $n^{th}$ C-UAV is determined as,

$$R_{mn}^{C-UAV}(t) = \frac{B_u}{M_n(t)} \log_2(1 + \Gamma_{nm}^{C-UAV}) \qquad (4)$$

here $M_n(t)$ represents the no. of active UEs served by $n^{th}$ C-UAV at time instance $t$.

- **C-UAV to BS transmission via I-UAV:** The desired signal from $n^{th}$ C-UAV to BS with power $P_n^t$ for further computation via $p^{th}$ I-UAV is denoted as $\hat{y}_{C-UAV}^{np}$. The jamming signal $x_J$ is transmitted to the BS with power $P_J$, undetectable by the C-UAV. Each IRS element reflects a combined signal to both the BS and the eavesdropper. The $p^{th}$ IRS's reflection coefficient matrix is $\mathbf{v_p} = (v_{1p}, \ldots, v_{Lp})^T$, where $v_{ip} = e^{j\theta_{ip}}$ and

---

t, illustrated in Fig. 1, comprises a distributed mobile user network with $M$ UEs positioned at $s_m^{UE}(t)$, $N$ C-UAVs deployed at $s_n^{C-UAV}(t)$ to provide Mobile Edge Computing (MEC) services, and a single Base Station (BS) at the origin. The system operates in the presence of a potential Jammer (J) at $s^J(t)$ and a single-antenna Eavesdropper (E) at $s^E(t)$. The exact locations of eavesdropper and jammer remain unknown. Instead, UAVs employ aerial photography target detection techniques to estimate their approximate positions, subsequently sharing this information with the BS. In our assumptions, each eavesdropper's and jammer's estimated regions, are denoted as $s^E(t)$ and $s^J(t)$, possesses a known radius $\epsilon_E$ and $\epsilon_J$ as discerned by the UAV, where $\epsilon_E \geq \|s^{\hat{E}}(t) - s^E(t)\|$ and $\epsilon_J \geq \|s^{\hat{J}}(t) - s^J(t)\|$, here $s^{\hat{E}}(t)$ and $s^{\hat{J}}(t)$ denotes the centroid of E's and J's potential region respectively. To overcome challenges in urban communication, we deploy $P$ IRS aided UAVs (I-UAVs). These UAVs are strategically positioned at $s_p^{I-UAV}(t)$ to create virtual Line-of-Sight (LoS) communication, addressing issues like data interception. Equipped with $L$ reflecting elements, the IRS, controlled by a microchip, manages communication between C-UAVs and the base station (BS). Trajectory planning for C-UAVs is crucial, considering time-varying data generation and limited energy. Coordination prevents collisions, and offloading starts when UAVs reach optimal locations. Each UE manages computation tasks periodically, with $m^{th}$ UE handling assignments represented by $W_m = [C_m, \lambda_m, D_m]$, here $D_m$ represents task data size, $C_m$ is CPU cycle count, and $\lambda_m$ is task arrival rate.. Due to limited capacity, UEs offload tasks to C-UAVs, which, constrained by size and weight, provide partial computation. C-UAVs then offload remaining tasks to the BS via I-UAVs for further processing. Our system, using Orthogonal Frequency Division Multiple Access (OFDMA), ensures simultaneous execution of tasks by multiple users, maintaining data freshness and security.

$\theta_{ip} \in [0, 2\pi]$ with $|v_{ip}| = 1$ for all $i$. Potential collaboration between the jammer and eavesdropper nullifies the eavesdropper's reception of the jamming signal. Multiple reflections by the IRS are negligible due to significant path loss. Consequently, the received signals at the BS and the eavesdropper by $n^{th}$ C-UAV are expressed as

$$y_n^{BS} = \sum_{p=1}^{P} \frac{h_{p0}^{IB}}{\sqrt{L_{p0}^{IB}}} \cdot diag(\mathbf{v_p}) \cdot \frac{h_{np}^{CI}}{\sqrt{L_{np}^{CI}}} \cdot \hat{y}_{C-UAV}^{np}$$
$$+ \sum_{p=1}^{P} \frac{h_{p0}^{IB}}{\sqrt{L_{p0}^{IB}}} \cdot diag(\mathbf{v_p}) \cdot \frac{h_{0p}^{JI T}}{\sqrt{L_{0p}^{JI}}} \cdot x_J$$
$$+ \frac{h^{JB}}{\sqrt{L^{JB}}} \cdot x_J + n_{BS} \quad (5)$$

$$y_n^E = \sum_{p=1}^{P} \frac{h_{p0}^{IE}}{\sqrt{L_{IE}^{p0}}} \cdot diag(\mathbf{v_p}) \cdot \frac{h_{np}^{CI}}{\sqrt{L_{np}^{CI}}} \cdot \hat{y}_{C-UAV}^{np} + n_E \quad (6)$$

here $n_{BS} \sim N(0, \sigma_{BS}^2)$ and $n_{BS} \sim N(0, \sigma_{BS}^2)$. The instantaneous received SNR at BS and eavesdropper receiving signal from $n^{th}$ C-UAV via $p^{th}$ I-UAV will be,

$$\Gamma_{np}^{BS}(t, v_p) = \frac{(\tilde{h}_{p0}^{IB} \cdot diag(v_p) \cdot \tilde{h}_{np}^{CI})^2 \cdot P_n^t}{(\sigma_{BS})^2 + \tilde{h}^{JB 2} \cdot P_J + \Phi_1} \quad (7)$$

$$\Gamma_{np}^{E}(t, v_p) = \frac{(\tilde{h}_{p0}^{IE} \cdot diag(v_p) \cdot \tilde{h}_{np}^{CI})^2 \cdot P_n^t}{(\sigma_E)^2 + \Phi_2} \quad (8)$$

where $\tilde{h}$ is path loss incorporated channel coefficient, $\Phi_1 = \sum \sum [(\tilde{h}_{p0}^{IB}).diag(\mathbf{v_p}).\tilde{h}_{np}^{CI}]^2 \cdot P_n^t + \sum [(\tilde{h}_{p0}^{IB}).diag(\mathbf{v_p}).\tilde{h}_{0p}^{JI}]^2 \cdot P_J$, and $\Phi_2 = \sum \sum [\tilde{h}_{p0}^{IE} \cdot diag(v_p) \cdot \tilde{h}_{np}^{CI 2} \cdot P_n^t]$. The instantaneous data rate between $n^{th}$ C-UAV and BS having bandwidth $B_{BS}$ via $p^{th}$ I-UAV given as,

$$R_{np}^{BS}(t, v_p) = \frac{B_{BS}}{N} \log_2(1 + \Gamma_{np}^{BS}(t, v_p)) \quad (9)$$

$$R_{np}^{E}(t, v_p) = \log_2(1 + \Gamma_{np}^{E}(t, v_p)) \quad (10)$$

Secrecy rate of $n^{th}$ C-UAV data at BS will be given as,

$$R_{sec,n}(t, v_p) = \sum_{p=1}^{P} \max\{0, [R_{np}^{BS}(t, v_p) - R_{np}^{E}(t, v_p)]\} \quad (11)$$

### B. Task Offloading Protocol

The overall task offloading process from UEs to BS has three phases. In the first phase, UEs communicate with C-UAVs. Second phase is subdivided into two parts, i.e. computation of tasks at C-UAVs, and simultaneously transmission from C-UAVs to BS via I-UAVs. Finally, computation at BS.

*1) Transmission from UE to C-UAV (U2C):* The time delay, and energy consumption associated with transmission between $m^{th}$ UE and $n^{th}$ C-UAV will be expressed as,

$$T_{mn}^{U2C} = \frac{D_m}{R_{mn}^{C-UAV}(t)} \quad (12)$$

$$E_{mn}^{U2C} = h_{mn}^{UC 2} \cdot P_m \cdot T_{mn}^{U2C} \quad (13)$$

*2) Computation at C-UAV:* Once C-UAVs recieve all data from the UEs, each C-UAV decides the amount of task that can be computed locally. The proportion of task of $m^{th}$ UE computed at $n^{th}$ C-UAV is given as $\beta_{m0}^n \in [0, 1]$, and the task executed at base station transmitted via $p^{th}$ I-UAV is given by $\beta_{mp}^n \in [0, 1]$, such that

$$\beta_{m0}^n + \sum_{k=1}^{K} \beta_{mk}^n = 1, \forall n \quad (14)$$

The computation delay at $n^{th}$ C-UAV while handling the task of $m^{th}$ UE is given by

$$T_{mn}^{C-UAV}(t) = \frac{\beta_{m0}^n D_m C_m}{f_{mn}(t)} \quad (15)$$

$f_{mn}(t) = \frac{F_u}{M_n(t)}$ represents the computational resources of *C-UAV n* allocated to $m^{th}$ UE, where $F_u$ is the computational resource of each C-UAV allocated equally to every UE.

Next, by taking into account the computation time and power consumption. The energy consumed by $n^{th}$ C-UAV while handling the task of $m^{th}$ UE.

$$E_{mn}^{C-UAV}(t) = \kappa[f_{mn}(t)]^3 T_{mn}^{C-UAV}(t) \quad (16)$$

where $\kappa$ stands for effective switched capacitance [9].

*3) Transmission from C-UAV to BS via I-UAV (C2I):* Taking into account that some tasks are offloaded by a particular $n^{th}$ C-UAV for further computing to BS via $p^{th}$ I-UAV. The time delay and energy consumption for this transmission will be given as

$$T_{np}^{C2I}(t, v_p) = \frac{\beta_{mp}^n \cdot D_m}{R_{np}^{BS}(t, v_p)} \quad (17)$$

$$E_{np}^{C2I}(t, v_p) = T_{np}^{C2I}(t, v_p)[(h_{p0}^{IB} \cdot diag(v_p) \cdot h_{np}^{CI})^2 \cdot P_n^t] \quad (18)$$

*4) Computation at BS:* After task data is offloaded by C-UAVs to BS, the BS starts processing the computation task. The computation delay at the BS is determined in terms of task ratio $\beta_{mp}^n$ as

$$T_{mnp}^{BS}(t) = \frac{\beta_{mp}^n D_m C_m}{f_n(t)} \quad (19)$$

where $f_n(t) = \frac{F_{BS}}{N}$ is the computational resource allocated to $n^{th}$ C-UAV at BS. $F_{BS}$ is the resource available at BS which is divided equally among all C-UAV.

It is assumed that each C-UAV has distinct computation and communication units. As a result, computations can be executed concurrently with the transmission of tasks. So, the computational time of $n^{th}$ C-UAV will be,

$$\tau_{t,comp}^n = \sum_{m=1}^{M} [T_{mn}^{U2C} + \max\{T_{mn}^{C-UAV}(t), T_{mnp}^{C2I}(t) + T_{mnp}^{BS}(t)\}] \quad (20)$$

## C. UAV movement

A task is divided into $T$ timeslots of length $\tau$, denoted as $\{0, 1, \ldots, t, \ldots, T-1\}$, where $t$ is the index of the current timeslot. Initially, all UAVs (C-UAVs and I-UAVs) are deployed at the origin. In each timeslot, an $m^{th}$ UE either generates a data packet $D_m$ or remains inactive. Data leaves the queue only when C-UAVs approach and collect it. Each C-UAV $n$ spends $\tau_{t,move}^n$ moving in direction $\rho_t^n \in (0, 2\pi)$ at a fixed speed. During remaining time $\tau_{t,comp}^n$, C-UAVs collect, compute, and transmit data. Similarly, I-UAVs take $\tau_{t,move}^p$ to reach a location and then hover. UEs transmit complete packets each time. Each UAV has limited energy reserve $E_{\max}$; the task fails if any UAV runs out of energy. The energy consumption of $i^{th}$ UAV due to flying and hovering in timeslot $t$ will be,

$$E_i^{UAV}(t) = \alpha_{\text{move}} \tau_{t,move}^i + \alpha_{\text{hover}}[\tau - \tau_{t,move}^i] \quad (21)$$

where $\alpha_{\text{move}}$, and $\alpha_{\text{hover}}$ are energy consumption coefficient while UAV is moving and hovering respectively. These coefficients are calculated using [10].

$$\alpha = c_1 \left(1 + \frac{3v_{\text{uav}}^2}{v_{\text{tip}}^2}\right) + c_2 \left(\sqrt{1 + \frac{v_{\text{uav}}^4}{4v_0^4} - \frac{v_{\text{uav}}^2}{2v_0^2}}\right) + \frac{1}{2}c_3 v_{\text{uav}}^3 \quad (22)$$

where $v_{\text{uav}} = v_{\text{move}}$ or $v_{\text{uav}} = 0$ for $\alpha_{\text{move}}$ and $\alpha_{\text{hover}}$. Constants $c_1$, $c_2$, $c_3$ depend on power, rotors, and air density. $v_{\text{tip}}$ is tip speed and $v_0$ denotes average velocity induced by the rotor.

## III. PROBLEM FORMULATION AND TRANSFORMATION

Initially, we present the metrics used in this paper. As previously highlighted, our specific focus revolves around data freshness and secure communication in presence of $J$ and $E$.

- **Threshold AoI Violation:** Let $z_m(x) \in \mathbb{R}$ be a random process in $x \in [0, X]$ represents the generation time of the oldest data at $m^{th}$ UE, where $X = T \cdot \tau$. The term $(x - z_m(x))$ signifies waiting time before data is collected. Let $G_m$ denote timeslot set where AoI threshold is violated at $m^{th}$ UE, i.e., $G_m = \{x | x \in [0, X], x - z_m(x) \geq AoI_{th}\}$. The violation ratio $\chi$ will be:

$$\chi = \frac{1}{M} \sum_{m=1}^{M} \int_0^X 1_{V_p}(x)\, dx, \quad (23)$$

where $1_{V_p}(\cdot)$ is the indicator function.

- **AoI Penalty:** A transformation applied to AoI, where $(\gamma)$ is the soft-constrained penalty function. In this context, the expression will be given as:

$$Q = \frac{1}{M} \sum_{m=1}^{M} w_m \int_0^X f(x - z_m(x))\, dx, \quad (24)$$

where weight $w_m$ is given to higher priority UEs and:

$$f(x - z_m(x)) = \begin{cases} \tilde{\gamma}(x - z_m(x)), & \text{if } x \in G_m \\ x - m_p(x), & \text{otherwise} \end{cases} \quad (25)$$

here $\tilde{\gamma}(x) = x + e^x$, introduce an exponential penalty gives more priority to UEs which violates AoI threshold.

## A. Problem Statement

Our goal is to minimize the threshold AoI violation and AoI Penalty for data freshness, and maximize achievable secrecy rate for secure communication under energy and motion constraint. This will be expressed as,

$$\min_{m,p,\beta,P_n^t,s_{i,UAV}(t),v_p} Q + \chi - \sum_{n=1}^{N} R_{sec,n}(t, v_p) \quad (26)$$

s.t. C1: $\sum_{m=1}^{M} E_{mn}^{U2C} + E_{mn}^{C-UAV} + \sum_{p=1}^{P} E_{np}^{C2I}$
$$+ E_n^{UAV} \leq E_{max}, \forall n \in [1, N]$$

C2: $E_p^{UAV} \leq E_{max}, \forall p \in [1, P]$

C3: $\sum_{p=1}^{P} R_{np}^E(t, v_p) \leq R_{th}, \forall n \in [1, N]$

C4: $\| s_{n1}^i(t) - s_{n2}^i(t) \| \geq D_{min}$
$$\forall n1, n2, n1 \neq n2, i \in \{C - UAV, I - UAV\}$$

C5: $C_{n1_{max}} + C_{n2_{max}} \leq \| s_{n1}^i(t) - s_{n2}^i(t) \|$
$$\forall n1, n2, n1 \neq n2, i \in \{C - UAV, I - UAV\}$$

In our formulation, constraints C1 and C2 impose maximum energy limits on both the C-UAV and I-UAV. Subsequently, C3 sets an upper boundary on the eavesdropper's data transmission rate. Constraints C4 and C5 delineate collision and overlapping restrictions, where $D_{\min}$ represents the minimum distance between two UAVs, and $C_{n1_{\max}}, C_{n2_{\max}}$ denote the coverage radius of the two nearest UAVs. Solving the intricate NP-hard optimization problem (26) involves numerous unknowns, including the UE's location and channel conditions. UAV mobility introduces dynamism, adding complexity. Traditional optimization methods struggle with the multitude of possible solutions. To address this, the next section delves into a DRL approach, aiming to formulate a near-optimal policy with minimal environmental information.

## B. Markov Decision Problem (MDP) Formulation

To address the problem outlined above, we initially represent it as a decentralized observable MDP, denoted as $M = \langle U, O, A, R, \Omega, \gamma \rangle$, where $U$, $\Omega$ and $\gamma$ denote combine UAV set containing C-UAV and I-UAV, transition probabilities, and the discounted factor, respectively.

- **Observation Space** $(O)$: The observation space $O$ is defined as $\{o_t\}$. Each UAV maintains its local observation $o_{ut}$ within a fixed sensing range, expressed as a disjoint union of $o_{ut}(\text{env})$ and $o_{ut}(\text{UAV})$. The former encompasses location, remaining data amount, and data generation time for all UEs within the sensing range, along with the current region of the Jammer and Eavesdropper, and jamming power $(P_J)$. The latter includes the current position and remaining energy for all UAVs during training.

- **Action Set** $(A)$: The action space $A \equiv \{a_t, v_p\}$. For each UAV, $a_{ut} = (\rho_t^u, l_t^u)$, where $\rho_t^u$ is the angle controlling the

direction of UAV movement, and $l_t^u$ is traveling distance, bounded by maximum distance $l_{\max}$, and $v_p$ is $p^{th}$ IRS's reflection coefficient matrix for an episode.

- **Reward Function** $(R)$**:** This system incorporates two types of reward functions: intrinsic and extrinsic.

  1) **Extrinsic Reward:** The environment provides the following reward for each UAV:

$$R_{\text{env},u}(t) = -\{Q + \chi - \sum_{n=1}^{N} R_{\text{sec},n}(t, v_p)\} \quad (27)$$

  2) **Intrinsic Reward:** This considers all penalties associated with UAV constraints:

$$R_{\text{int},u}(t) = -\{\eta_1 + \eta_2 + \eta_3\} \quad (28)$$

Here, $\eta_1, \eta_2, \eta_3$ represent constant penalties applied if constraints $C_1, C_2, C_3$ are violated, respectively.

In this context, the objective is to maximize the sum of overall reward $R_u(t) = R_{\text{int},u}(t) + R_{\text{env},u}(t)$.

## IV. GTR-DRL ALGORITHM

We present a novel decentralized Multi-Agent DRL (MADRL) approach for our system, integrating a Gated Transformer (GTr) for temporal modeling. Decentralized policies commonly encounter slow convergence issues. To address this challenge, we initially implement IMPALA principles [11] in a multi-agent context. Our proposed algorithm follows the *Decentralized Training and Collaborative Execution* (DTCE) paradigm, with asynchronous actors operating on separate GPUs to enhance training efficiency.

### A. Transformer-Enhanced Distributed Framework

The architecture incorporates GTr blocks, crucial for extracting temporal features from past trajectories. Each UAV's observation undergoes embedding via multi-layer perceptions (MLP), and temporal features are derived through multi-head attention (MHA) as shown in Fig. 2. Layer normalization is applied for model stabilization. The decentralized framework integrates policy and value networks for each UAV, addressing issues related to policy lagging and estimation variance.

V-trace target is employed to address high variance challenges in UAV optimization. Our approach utilizes $L_2$ loss for decentralized value network optimization and policy gradients for updating the decentralized policy network. Sharing parameters between networks contributes to improved training efficiency. The equation used for calculating the V-trace target for each UAV is given as,

$$V_{u,t}(\mathbf{o}_u) = V_{u,\phi}(\mathbf{o}_u) + \sum_{i=t}^{T-1} \gamma^{i-t} \left( \prod_{j=t}^{i-1} c_{u,j} \right) \rho_{u,i} \cdot \text{TD}_{u,i}$$

here $V_{u,t}(\mathbf{o}_u)$ represents the V-trace target for UAV $u$ at time $t$, computed based on the observation $\mathbf{o}_u$. $V_{u,\phi}(\mathbf{o}_u)$ denotes the value estimate derived directly from the current observation for UAV $u$. $c_{u,j}$ corresponds to the truncated importance sampling ratio at time step $j$. $\rho_{u,i}$ signifies the truncated importance
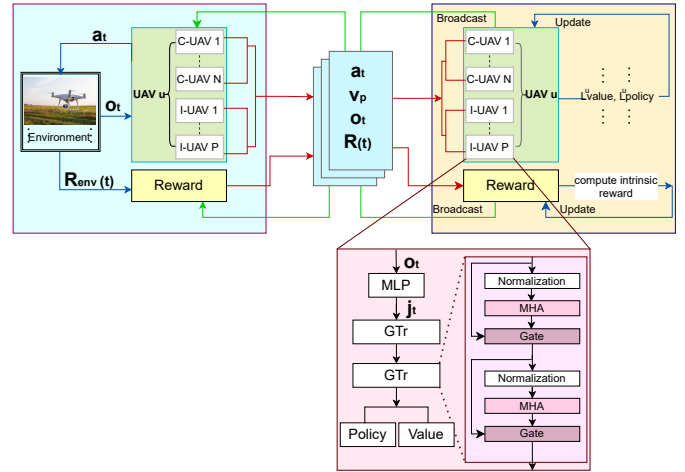


Fig. 2: Proposed Solution Network

sampling ratio at time step $i$. $\text{TD}_{u,i}$ stands for the one-step temporal difference target at time step $i$. The framework is mathematically represented with,

1) **Decentralized Policy Update:**

$$L_{\text{policy},u}(\theta) = E\left[\rho_{u,t} \log \pi_{u,\theta}(a_{u,t}|o_{u,t}) \right.$$
$$\left. (r_{u,t} + \gamma V_u(o_{u,t}) - V_{u,\phi}(o_{u,t}))\right] \quad (29)$$

2) **Decentralized Value Network Optimization:**

$$L_{\text{value},u}(\phi) = E\left[(V_u(o_{u,t}) - V_{u,\phi}(o_{u,t}))^2\right] \quad (30)$$

### B. Optimization

The framework consists of independent entities called actors, operating UAVs and IRS matrices, responsible for local information gathering. A central decision-making entity, the learner, utilizes collected experiences to update deep neural network (DNN) weights, enabling decentralized and collaborative learning.

The Actor in Algorithm 1, embodies the behavior of individual UAVs. Initialized with its unique policy network weights $(\theta_{\text{act}})$, episodic buffer, and other parameters, the Actor continuously interacts with the environment, gathering experiences that include UAV observations, selected actions, and environmental rewards. When the episodic buffer is full, the actor communicates with the learner by sending the experiences. If the learner broadcasts updated network weights $(\theta)$, the actor updates its policy network accordingly. The Learner, as depicted in Algorithm 2, functions as the central learning entity. It initializes DNN weights $(\theta$, in a learning loop, aggregates experiences from multiple actors. The learner computes value and policy loss functions $(L_u^{\text{value}}(\phi)$ and $L_u^{\text{policy}}(\theta))$ based on specific equations for each UAV. Using gradient descent methods, it minimizes the weighted sum of all losses $(L_{\text{total}})$. To ensure synchronization, the learner periodically broadcasts the updated network weights $(\theta)$ to all actors. This framework adapts to diverse and dynamic scenarios through decentralized learning and collaboration.

**Algorithm 1** Actor

1: **procedure** ACTOR($\theta_{\text{act}}$, episodic buffer)
2:     **while** learner updates **do**
3:         Clear episodic buffer
4:         **while** episodic buffer is not full **do**
5:             Get UAVs' observation $o_t$ and select actions $a_t, v_p$ from policy $\pi_{\theta_{\text{act}}}$
6:             Interact with the environment and get reward $R_{u,env}(t)$
7:             Compute total rewards $R_u(t)$ by Equations (26) and (27) and store experiences in the buffer
8:         **end while**
9:         Send full episodic buffer to learner
10:        **if** received broadcast weights $\theta$ **then**
11:           Update network $\theta_{\text{act}} \leftarrow \theta$
12:        **end if**
13:     **end while**
14: **end procedure**

---

**Algorithm 2** Learner

1: **procedure** LEARNER($\theta$)
2:     Initialize network weights $(\theta)$
3:     **while** learner updates **do**
4:         Get experiences from different actors
5:         Compute $L_u^{\text{value}}(\phi), L_u^{\text{policy}}(\theta)$ for each UAV by Equations (28) and (29)
6:         Minimize weighted sum of all losses $L_{\text{total}}$ by gradient descent methods
7:         **if** updated times mod broadcast interval $= 0$ **then**
8:            Send network weights $\theta$ to all actors
9:         **end if**
10:     **end while**
11: **end procedure**



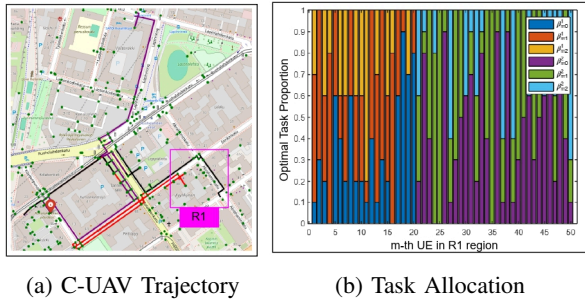(a) C-UAV Trajectory        (b) Task Allocation

Fig. 3: C-UAV Trajectory and Task Allocation

## V. RESULTS AND DISCUSSION

In this section, we present comprehensive simulation results, including the trajectory of C-UAVs, the optimized task allocation between C-UAVs and BS via I-UAVs within our proposed learning-oriented approach. Additionally, we offer insights into key data freshness metrics, such as threshold AoI violation and AoI Penalty, alongside the average secrecy rate achieved by our learning-based scheme. These results are compared against five benchmark schemes, namely DDPG-EWSA[12], IMPALA[11], VDN [13], QMIX[12], and random approach. In our simulation, the system operates within a designated area of $1500\,\text{m} \times 1500\,\text{m}$, with all UAVs initially positioned at the origin. UE distribution follows a normal distribution across the region, with UE positions generated at the outset and remaining constant throughout each time interval. Learning system parameters are set with a mini-batch size of 256, a replay memory size of $2 \times 10^5$, $1 \times 10^3$ training episodes, a learning rate of 0.005, a discount factor of 0.95, and a soft update rate of 100. Neural networks in the learning system comprise two hidden layers with ReLU activation functions, utilizing the Adam optimizer. City data, including the locations and configurations of tall buildings, is annotated using Google Maps and OpenStreetMap to ensure UAV collision avoidance.

Fig. 3 illustrates the trajectories of three C-UAVs collaborating with a total of six I-UAVs for security support. The red icon indicates the starting point of the C-UAVs. Additionally, the figure includes a task allocation plot for region $R1$, where 50 UEs are situated. In this region, two C-UAVs and two I-UAVs are tasked with specific assignments. For instance, the task allocation for the 50th UE showcases the collaborative deployment approach: 10% of the task is handled by C-UAV 2, 20% is transmitted to the BS via I-UAV 1, and the remaining 70% is efficiently transferred by I-UAV 2.

Fig. 4 illustrates the simulation results in the presence of jammers and eavesdroppers, showcasing the impact on the AoI penalty, threshold AoI violation, and secrecy rate as the number of C-UAVs increases. In this scenario, we maintain a fixed AoI threshold ($AoI_{th} = 100$) and deploy 6 I-UAVs. Notably, the results reveal improved data freshness metrics. Fig. 4(a) shows decrease in AoI Penalty which means the delay in data collection drops, and in Fig. 4(b) we see that ratio of UEs which exceed AoI threshold significantly reduced, but average secrecy rate also goes down as shown in Fig. 4(c), particularly in denser C-UAV network. The findings suggest a crucial trade-off between data freshness and security the system while determining the optimal number of C-UAVs.

Fig. 5, presents performance metrics concerning changes in the number of I-UAVs while maintaining a constant presence of 3 C-UAVs. Notably, there is no significant change in AOI-based metrics with varying numbers of I-UAVs depict in Fig. 5(a) and (b). However, average secrecy rate shown in Fig. 5(c), increases with increase in I-UAVs, providing enhanced support. But beyond a certain point, the average secrecy rate begins to decline. This decline is due to increase in unknown channels interacting with third-party nodes. This shows trade-off between increasing I-UAV support for security and mitigating potential risks associated with additional unknown channels. Deploying more UAVs amplifies the optimization challenge in an exponentially expanding solution space. Our algorithm, GTr-DRL, outshines baseline methods: DDPG-EWSA, employing DDPG with Ornstein-Uhlenbeck noise, and VDN and QMIX with $\epsilon$-greedy policies proved insufficient in this dynamic environment. IMPALA, while surpassing DDPG-EWSA, faced non-stationary problems due
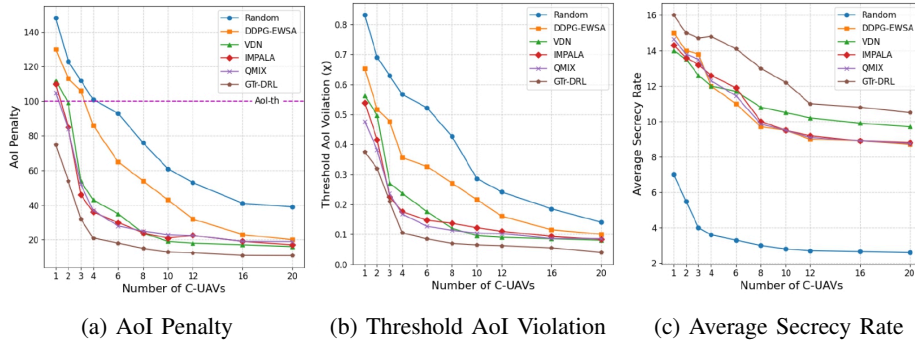
(a) AoI Penalty     (b) Threshold AoI Violation     (c) Average Secrecy Rate

Fig. 4: Performance parameters vs number of C-UAVs



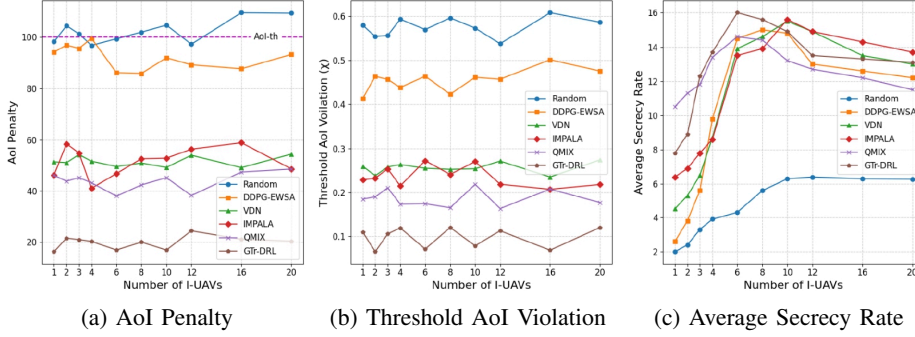(a) AoI Penalty     (b) Threshold AoI Violation     (c) Average Secrecy Rate

Fig. 5: Performance parameters vs number of I-UAVs

to centralized paradigms as neglecting specific UAVs' behavior by just relying on global joint reward. Despite enhancements, QMIX encountered hurdles from $\epsilon$-greedy exploration and lack of temporal modeling. This underscores GTr-DRL's effectiveness in multi-UAV netwrok optimisation.

## VI. CONCLUSION

In this paper, an efficient approach with IRS-assisted AoI-aware bi-layer multi-UAV system is designed to optimize task offloading and enhance overall network performance. The framework introduces exponential AoI metrics and prioritizes the maximization of secrecy rates, addressing critical challenges in UAV-assisted networks. By using multi-agent DRL, our approach efficiently allocates tasks, minimizing AoI, and ensuring robust security. The study highlights the trade-off between AoI and secrecy rate, emphasizing the need for balancing information delay and data confidentiality.

## REFERENCES

[1] B. Choudhury, V. K. Shah, A. Ferdowsi, J. H. Reed, and Y. T. Hou, "AoI-minimizing scheduling in UAV-relayed IoT networks," in *2021 IEEE 18th International Conference on Mobile Ad Hoc and Smart Systems (MASS)*. IEEE, 2021, pp. 117–126.

[2] O. Rahimi and A. Shafieinejad, "Minimizing age of information in multi-UAV-assisted IoT networks: a graph theoretical approach," *Wireless Networks*, pp. 1–23, 2023.

[3] M. Sun and et. al, "AoI-Energy-Aware UAV-Assisted Data Collection for IoT Networks: A Deep Reinforcement Learning Method," *IEEE Internet of Things Journal*, vol. 8, no. 24, pp. 17 275–17 289, 2021.

[4] W. Wei, X. Pang, J. Tang, N. Zhao, X. Wang, and A. Nallanathan, "Secure transmission design for aerial irs assisted wireless networks," *IEEE Transactions on Communications*, 2023.

[5] H. Shakhatreh, A. Sawalmeh, A. H. Alenezi, S. Abdel-Razeq, and A. Al-Fuqaha, "Mobile-IRS assisted next generation UAV communication networks," *Computer Communications*, 2023.

[6] W. Jiang, B. Ai, M. Li, W. Wu, and X. Shen, "Average age-of-information minimization in aerial irs-assisted data delivery," *IEEE Internet of Things Journal*, vol. 10, no. 17, pp. 15 133–15 146, 2023.

[7] N. Beaulieu and C. Cheng, "Efficient Nakagami-m fading channel Simulation," *IEEE Transactions on Vehicular Technology*, vol. 54, no. 2, pp. 413–424, 2005.

[8] X. Yu and et. al, "Robust and secure wireless communications via intelligent reflecting surfaces," *IEEE Journal on Selected Areas in Communications*, vol. 38, no. 11, pp. 2637–2652, 2020.

[9] J. Xiong, H. Guo, and J. Liu, "Task offloading in uav-aided edge computing: Bit allocation and trajectory optimization," *IEEE Communications Letters*, vol. 23, no. 3, pp. 538–541, 2019.

[10] M. Mozaffari, W. Saad, M. Bennis, Y.-H. Nam, and M. Debbah, "A tutorial on UAVs for wireless networks: Applications, challenges, and open problems," *IEEE communications surveys & tutorials*, vol. 21, no. 3, pp. 2334–2360, 2019.

[11] L. Espeholt, H. Soyer, R. Munos, K. Simonyan, Y. Mnih, V. Firoiu, T. Harley, I. Dunning *et al.*, "Impala: Scalable distributed deep-RL with importance weighted actor-learner architectures," in *International conference on machine learning*. PMLR, 2018, pp. 1407–1416.

[12] T. Rashid, M. Samvelyan, C. S. De Witt, G. Farquhar, J. Foerster, and S. Whiteson, "Monotonic value function factorisation for deep multi-agent reinforcement learning," *The Journal of Machine Learning Research*, vol. 21, no. 1, pp. 7234–7284, 2020.

[13] P. Sunehag, G. Lever, A. Gruslys, W. M. Czarnecki, V. Zambaldi, M. Jaderberg, M. Lanctot, N. Sonnerat, J. Z. Leibo, K. Tuyls *et al.*, "Value-decomposition networks for cooperative multi-agent learning," *arXiv preprint arXiv:1706.05296*, 2017.