

Resource Efficient Actuation in UAV-aided Sensor-Actuator Networks

Original

Resource Efficient Actuation in UAV-aided Sensor-Actuator Networks / Goel, A., De, S., Chiasserini, C., Casetti, C.E.. - In: IEEE COMMUNICATIONS LETTERS. - ISSN 1089-7798. - 29:10(2025), pp. 2223-2227. [10.1109/LCOMM.2025.3590648]

Availability:

This version is available at: 11583/3001812 since: 2025-10-12T02:48:21Z

Publisher:

IEEE

Published

DOI:10.1109/LCOMM.2025.3590648

Terms of use:

This article is made available under terms and conditions as specified in the corresponding bibliographic description in the repository

Publisher copyright

IEEE postprint/Author's Accepted Manuscript

©2025 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collecting works, for resale or lists, or reuse of any copyrighted component of this work in other works.

(Article begins on next page)

Resource-Efficient Actuation in UAV-Aided Sensor-Actuator Networks

Amit Goel¹, *Student Member, IEEE*, Swades De¹, *Senior Member, IEEE*,
Carla-Fabiana Chiasserini², *Fellow, IEEE*, Claudio E. Casetti³, *Senior Member, IEEE*

Abstract—Integrating uncrewed aerial vehicles (UAVs) into wireless sensor-actuator networks (WSANs) offers flexibility and improved system performance. However, imprecise localization and limited battery capacity constrain the UAV operation. This paper explores the use of battery swap station (BSS)-assisted UAV to facilitate timely actuation in WSANs while addressing the above limitations. The UAV collects data via backscatter communication from the sensor nodes and delivers to the energy-constrained actuator nodes along with the required energy. Incorporating UAV location uncertainty and Nakagami- m wireless channel fading, closed-form expressions are derived for the ergodic capacity in backscatter communication and the expected energy harvesting rate. To minimize the maximum delay in actuation, an optimization problem is formulated. To reduce complexity, the problem is transformed into an equivalent node visit sequence optimization and solved using sequential deep reinforcement learning (SDRL). We verify the accuracy of our analysis through Monte Carlo simulations. Our results show that the proposed SDRL-based strategy consistently offers reduced actuation delay with a significantly small computation overhead.

Index Terms—UAV, wireless sensor-actuator network, backscatter communication, wireless energy transfer, DRL

I. INTRODUCTION

Recent advances in wireless communication have enabled distributed sensing and actuation through wireless sensor and actuator networks (WSANs) [1]. However, these networks, which consist of wirelessly connected sensor and actuator nodes, are constrained by communication range, energy, and infrastructure costs. Conventional nodes with limited battery life need frequent replacements, which is infeasible in hazardous deployments. To overcome these issues, emerging technologies such as backscatter communication (BSC) and wireless energy transfer (WET) have gained attention [2], [3].

In WSANs, a key objective is to deliver sensor data to actuators for timely action while ensuring sufficient energy availability. The Age of Information (AoI) metric quantifies the freshness of sensor updates [4]. However, minimizing AoI alone does not ensure timely data-driven actuation. While [5] adds data uncertainty to AoI, it does not fully address end-to-end responsiveness. In [6], actuation delay is modeled

assuming a direct sensor-actuator communication link, which limits its utility in remote deployments. To address this, [7] proposes using a static controller to relay information. However, a static controller may not be sufficient for communication and actuation operations, as these power-intensive tasks require the controller to be close to the sensor and actuator nodes. To this end, UAVs can be employed in BSC- and WET-based WSANs. While UAV-assisted communication and WET have been individually studied, their joint application with monostatic backscatter in sensor data relaying under practical constraints such as limited UAV energy, imprecise localization, and low WET efficiency has received limited attention. Although approaches like establishing BSSs [8], [9] can extend UAV service time, optimizing UAV-aided WSAN operation under such constraints requires in-depth analysis.

To address the above challenges, this work contributes as follows: (i) Using stochastic geometry tools, we derive closed-form expressions for BSC ergodic capacity \bar{r}_{BSC} and the expected energy harvesting rate \bar{P}_{EH} , incorporating Nakagami- m channel fading and UAV location uncertainty. (ii) Based on \bar{r}_{BSC} and \bar{P}_{EH} and WSAN requirements, we propose a novel UAV-assisted framework to minimize the actuation delay. (iii) To achieve the objective, we reformulate the problem into an equivalent form that optimizes the node visiting sequence. (iv) Given the high complexity due to the constrained combinatorial nature of the problem, a sequential deep reinforcement learning (SDRL) approach is employed for solving it. (v) Our numerical results demonstrate significant performance benefits with the proposed method compared to the benchmark schemes. While this study focuses on a single UAV, future multi-UAV extensions [10] offer promising scalability. We remark that the proposed system is practical for scenarios like avalanche-prone areas, enabling timely, sensor-driven warning actuation for road safety.

II. SYSTEM MODEL

Consider a UAV-aided WSAN, having N sensor-actuator pairs, a UAV, and a BSS, deployed in a circular region of radius R_{max} (Fig. 1). The nodes are assumed to be located far apart; thus, the UAV can serve only one node at a time. The UAV visits each sensor before its corresponding actuator once per round while accessing the BSS as required. It follows predefined inter-node paths at a speed V_{UAV} and altitude h_{UAV} , enabling low control complexity and energy-aware mobility. UAV power consumptions for movement P_{mov} and hovering P_{hov} are modeled as in [11]. The sensor and actuator nodes

¹Dept. of Electrical Engineering and Bharti School of Telecommunication Technology and Management, IIT Delhi, India (bsz208009@iitd.ac.in; swadesd@ee.iitd.ac.in). ²Dept. of Electronics and Telecommunications, Politecnico di Torino, Italy (chiasserini@polito.it). ³Dept. of Control and Computer Engineering, Politecnico di Torino, Italy (claudio.casetti@polito.it).

This work was supported in part by the Dept. of Science and Technology India-Italy Research Mobility under Grant INT/Italy/P-34/2022; in part by the Tehri Hydroelectric Development Corporation under Grant THDC/RKSH/R&D/F-2076/1036; and in part by the Science and Engineering Board under Grant CRG/2023/005421.

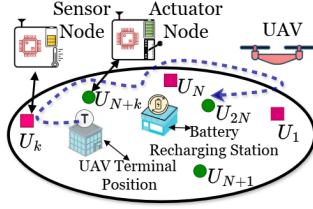


Figure 1 UAV and BSS aided wireless sensor actuator network.

have communication and energy harvesting modules; sensors include sensing units, whereas actuators contain mechanical actuation modules. The combined set of nodes is denoted by $\mathcal{U}=\{1, \dots, 2N\}$, where index $k \in \{1, \dots, N\}$ refers to the k -th sensor and $k+N \in \{N+1, \dots, 2N\}$ to its corresponding actuator. Nodes are assumed static with known positions, reflecting real-world scenarios where coordinates are fixed during deployment or obtained via localization methods, without loss of generality. The coordinates of n -th node is denoted as $w_n = \{x_n, y_n, z_n\}$. For convenience, indices '0' and '2N+1' represent BSS and UAV terminal station with respective coordinates $\{x_0, y_0, z_0\}$ and, $\{x_{2N+1}, y_{2N+1}, z_{2N+1}\}$. The 2D Euclidean distance between any two indices is denoted as $d_{i,j}$.

The UAV collects the sensor data using monostatic BSC by hovering over it and transmitting an unmodulated signal u_1 toward the sensor node. The received signal at the sensor node is $y_S = h_1 \sqrt{P_{tx}} u_1 + n_1$. The sensor node reflects back the received signal after modulating it with the information signal u_2 . The signal received by the UAV is [12]: $y_U = h_2 h_1 \sqrt{P_{tx}} u_1 u_2 + n_1 + n_2$, where P_{tx} is the transmit power of the UAV, and n_1 and n_2 denote additive white Gaussian noise (AWGN). There can be some imperfection in hovering position due to UAV localization error, which are discussed in Section III. **The LoS-dominant UAV channel is modeled by the Nakagami- m distribution, owing to its versatile fading representation and analytical tractability.** The channel gain is $h_i = \sqrt{G_1 G_2 (\lambda/4\pi)^2 d^{-\alpha_p} \tilde{h}_i} \forall i \in \{1, 2\}$. Here, \tilde{h}_i denotes the small-scale fading coefficients, α_p is the path loss exponent, λ is the transmission frequency, and G_1, G_2 , respectively, denote the UAV and node antenna gain. The communication data rates from the sensor node to UAV r_{S-U} and the UAV to actuator node r_{U-A} are expressed as $r_{S-U} = \log_2(1 + P_{tx} |h_1|^2 |h_2|^2 / \sigma_1^2 + \sigma_2^2)$, and $r_{U-A} = \log_2(1 + P_{tx} |h_1|^2 / \sigma_1^2)$. The harvested power P_{EH} , with $P_{tx} = P_{tx} |h_1|^2$, is modeled as [13]

$$P_{EH}(P_{tx}) = (aP_{tx} + b) \cdot \mathbf{1}_{[\varrho_1, \varrho_2)}(P_{tx}) + P_{EH}^{\max} \cdot \mathbf{1}_{[\varrho_2, \infty)}(P_{tx}) \quad (1)$$

where $\mathbf{1}_{[m, n)}(x) = 1$ if $m \leq x < n$, else 0, P_{EH}^{\max} is maximum harvest power, ϱ_1 and ϱ_2 are RF energy harvesting sensitivity and saturation threshold, a and b are shaping parameters.

III. STATISTICAL ANALYSIS

In this section, \bar{r}_{BSC} and \bar{P}_{EH} are derived based on the above channel model and UAV location uncertainty.

Distance distribution: While serving the n -th node, the UAV is meant to hover above w_n . Due to location error, its actual position is considered uniformly distributed in a ball

$B(w_n, \Delta_R)$. Thus, the UAV-to-node distance distribution is

$$f_d(d) = 2d / \sqrt{d^2 + h_{UAV}^2}; \quad h_{UAV} \leq d \leq \sqrt{h_{UAV}^2 + \Delta_R^2}. \quad (2)$$

BSC ergodic capacity: It is defined as the achievable average data rate, considering Nakagami- m wireless fading channel and the spatial distribution of the UAV-to-node distance as modeled in (2) and is expressed as $E_d [E_{h_1 h_2} [\log_2(1+\gamma)]]$.

Theorem 1. *The closed-form expression for BSC ergodic capacity is expressed in (4), where $C_3 = (\Gamma m_1 \Gamma m_2)^{-1}$, $C_4 = c_{41}/c_{42}$, $c_{41} = \alpha_1 \alpha_2 P_{tx} (G_1 G_2)^2 (\lambda/4\pi)^4$, and $c_{42} = m_1 m_2$.*

Proof. See Appendix A. ■

Expected energy harvesting rate: It is defined as the average amount of power that can be harvested by the node via WET process.

Theorem 2. *The closed-form expression for expected energy harvesting rate considering Gamma-distributed wireless fading link gain and spatial distribution as modeled in (2) is expressed in (5), where $K = P_{tx} E[d^{\alpha_p}] (\lambda/4\pi)^2 G_1$.*

Proof. See Appendix B. ■

IV. OPTIMIZATION FRAMEWORK AND CONSTRAINT

We now formulate the actuation delay minimization based on the node requirements and statistical measures from (4) and (5). Then we present the associated feasibility constraints.

The actuation delay minimization is formulated as a sequence prediction problem. After the UAV leaves the terminal, it travels to each sensor node to gather data and proceeds to the associated actuator node to transfer information and energy. This process continues until data are collected from all sensor nodes and the respective actuators are served, and the UAV recharges itself as required. UAV visiting sequence is denoted as $S = \{s_k\}_{k=1}^L$ where $s_k \in \{0, 1, \dots, 2N+1\}$ denotes the indices having one-to-one mapping with the nodes, BSS, and terminal station, L is the sequence length, which is at least $2(N+1)$, but could be longer since the UAV may need to visit the BSS multiple times. S needs to be optimized for reduced actuation delay. In finding the optimal S , the feasibility needs to be ensured, as the constraints can affect the viability of S .

1) *Feasible sequence:* We define a variable $\mathbf{I}(k)$ that indicates the index of the location visited by the UAV at k -th position in the sequence and is defined as $\mathbf{I}(k) = m$, s.t. $s_m = k$; $m \in \{1, \dots, L\}$, $k \in \{1, \dots, 2N\}$. While serving, each node should be visited only once, with an actuator visited only after its corresponding sensor is served. The constraints are represented as

$$\text{if } s_i = k \Rightarrow s_j \neq k, \forall j \neq i, \forall i, j \in \{1, \dots, L\}, k \in \{1, \dots, 2N\} \quad (6)$$

$$\mathbf{I}(k) < \mathbf{I}(N+k) \quad \forall k \in \{1, \dots, N\}. \quad (7)$$

Moreover, the total time the UAV spends on each node is required to be greater than or equal to the required service time for that node. The variable τ_j denotes the expected time required to serve the j -th node and is defined as

$$\bar{r}_{\text{BSC}} = \frac{C_3 h_{\text{UAV}}^2 \Delta_R^3}{2} H_{1,0,4,3;1,2}^{0,1:1,4:1,1} \left(\begin{matrix} (0, 2, 1) \\ - \end{matrix} \middle| \begin{matrix} (1 - m_1, 1) (1 - m_2, 1) (1, 1) (1, 1) \\ (1, 1) (1/2, 2) (0, 1) \end{matrix} \right) \left(\begin{matrix} (-1/2, 1) \\ (-1, 1) (-3/2, 1) \end{matrix} \right) C_4 h_{\text{UAV}}^{-4} \frac{\Delta_R^2}{h_{\text{UAV}}^2} \quad (4)$$

$$\bar{P}_{\text{EH}} = \frac{a\alpha_1 K}{m_1 \Gamma m_1} \left(\gamma \left(m_1 + 1, \frac{m_1 \varrho_2}{K \alpha_1} \right) - \gamma \left(m_1 + 1, \frac{m_1 \varrho_1}{K \alpha_1} \right) \right) + \frac{b}{\Gamma m_1} \left(\gamma \left(m_1, \frac{m_1 \varrho_2}{K \alpha_1} \right) - \gamma \left(m_1, \frac{m_1 \varrho_1}{K \alpha_1} \right) \right) + \frac{P_{\text{EH,max}}}{\Gamma m_1} \left(1 - \frac{1}{\Gamma m_1} \gamma \left(m_1, \frac{m_1 \varrho_2}{K \alpha_1} \right) \right) \quad (5)$$

$$\tau_j = \begin{cases} \frac{D_j^{\text{req}}}{\bar{r}_{\text{S-U},j}} \left(1 + \frac{P_c}{\bar{P}_{\text{EH}}} \right) & \forall j \in \{1, \dots, N\} \\ \frac{E_j^{\text{req}}}{\bar{P}_{\text{EH}}} + \frac{D_j^{\text{req}}}{\bar{r}_{\text{U-A},j}} \approx \frac{E_j^{\text{req}}}{\bar{P}_{\text{EH}}} & \forall j \in \{N+1, \dots, 2N\} \\ \tau_{\text{BSS}} & j = 0 \end{cases} \quad (8)$$

where D_j^{req} denotes the amount of data required to be transmitted, E_j^{req} denotes the required energy for actuation, P_c denotes the sensor node circuit power consumption while performing BSC [3], and τ_{BSS} denotes the time required for battery swapping. Thus, the total energy required to be spent by the UAV on the j -th node is expressed as $\mathcal{E}_j^{\text{ser}} \approx (P_{\text{Hov}} + P_{\text{tx}}) \tau_j \quad \forall j \in \{1, \dots, N\}$. The UAV transmits with P_{tx}^a for actuator nodes and P_{tx}^s for sensor nodes. Once a node has been served, $\mathcal{E}_j^{\text{ser}}$ is reset to null. Furthermore, throughout the process, it is necessary to monitor the on-board battery and recharge it when required. The variable $B_{\text{rem},j}$ denotes the remaining on-board battery capacity of UAV after it has completed traversing the sequence up to the j -th element, $\forall j \in \{0, 1, \dots, 2N+1\}$. Furthermore, to represent the relationship between two nodes that are visited in sequence, a binary variable $X_{ij}(n)$ is defined as $X_{ij}(n)=1$ if $s_n=i$ and $s_{n+1}=j$, where, $i, j \in \{0, \dots, 2N+1\}$ and $n \in \{0, \dots, L-1\}$. Given that $\mathcal{B}_{\text{rem},n}$ is known and $X_{ij}(n)=1$ where $n \in \{0, \dots, L-1\}$, then the battery status is updated as

$$\mathcal{B}_{\text{rem},n+1} = \begin{cases} \mathcal{B}_{\text{rem},n} - \mathcal{E}_{ij}^{\text{mov}} - \mathcal{E}_j^{\text{ser}} & j \in \{1, \dots, 2N\} \\ \mathcal{B}_{\text{full}} & j = 0 \\ \mathcal{B}_{\text{rem},n} - \mathcal{E}_{ij}^{\text{mov}} & j = 2N+1. \end{cases} \quad (9)$$

It is notable that the onboard battery status restricts the UAV movement. The constraints are represented as

$$X_{ij}(n) \geq 0 \text{ if } \begin{cases} \mathcal{B}_{\text{rem},n} \geq \mathcal{E}_{ij}^{\text{mov}} + \mathcal{E}_j^{\text{ser}} + \mathcal{E}_{j0}^{\text{mov}} & j \in \{1, \dots, 2N\} \\ \mathcal{B}_{\text{rem},n} \geq \mathcal{E}_{ij}^{\text{mov}} & j \in \{0, 2N+1\} \end{cases} \quad (10)$$

$$X_{i0}(n) = 1 \text{ if } \mathcal{B}_{\text{rem},n} - \mathcal{E}_{i,0} < \min \left\{ \mathcal{E}_{ij}^{\text{mov}} \right\} \forall n \in \{1, \dots, L-1\}, \quad (11)$$

$$i \in \{1, \dots, 2N\}, j \in \{1, \dots, 2N+1\} \setminus \{i\}$$

where (10) indicates that after serving the i -th node, the UAV can visit the j -th node only if the remaining battery capacity is sufficient for both serving the j -th node and returning to the BSS from that node. Moreover, (11) states that if the remaining battery capacity is not enough to serve any unserved nodes, the UAV should proceed to the BSS. The subtraction in (11) ensures that the UAV always has sufficient battery to reach the BSS. Furthermore, (12) highlights that the UAV in any feasible sequence starts and ends at the UAV terminal position.

$$\sum_i X_{i,2N+1,i}(0) = 1, \sum_i X_{i,2N+1}(L-1) = 1 \forall i \in \{0, \dots, 2N\}. \quad (12)$$

2) **Problem formulation:** The UAV performance is measured by the maximum delay of actuation (MDA), defined as the maximum time difference between the start of operation and when the actuator takes action. The overall optimization

problem considering this metric and the feasibility constraints is expressed as:

$$(\mathcal{P}1) : \min_{\mathbf{S}} \max_{\mathbf{k}} |t_{k+N} - t_0| \quad (13)$$

s.t. : (6), (7), (10) – (12)

where, t_0 is the time at which UAV starts the operations and t_{k+N} denotes the time at which actuator associated with k -th sensor takes action which is expressed as $t_j = \sum_{i=1}^{\mathbf{I}^{(j)}} (d_{s_i, s_{i+1}} / V_{\text{UAV}}) + \tau_i, \forall j \in \{N+1, \dots, 2N\}$. The problem $\mathcal{P}1$ is an NP-hard combinatorial optimization problem. Considering the UAV battery capacity sufficiently large and homogeneity among all the nodes, the transformed problem becomes an equivalent NP-hard problem as in [11]. Consequently, the problem $\mathcal{P}1$ is inherently NP-hard, rendering it difficult to address through conventional solution methods. Therefore, we propose a sequential deep reinforcement learning approach building on [14], [15] to effectively address this challenge. The details of the solution are outlined in the following section.

V. PROPOSED SDRL-BASED SOLUTION

The proposed solution utilizes a Markov Decision Process framework, where the UAV (agent) interacts with the dynamic environment by selecting actions that generate rewards while transitioning among states. Here, the UAV is responsible for making decisions and executing actions. A state S_l includes information about the last visited node and relevant parameters such as $\{w_n, \tau_n, \mathcal{E}_n^{\text{ser}}\}_{n=0}^{2N+1}$, and $B_{\text{rem},l}$. The action set consists of the available choices at a given state S_l , which includes unvisited node indices, BSS, and the UAV terminal station, all subject to the constraints defined in $\mathcal{P}1$. The action a_l represents the index of the location the UAV visits. The reward is calculated at the end of the episode and is defined as the negative of the maximum actuation delay, as described in $\mathcal{P}1$.

1) **Sequential neural network architecture:** To map the current state to a probability distribution of possible actions, a sequential neural network-based policy is used at the UAV. Primarily, sequential neural network architecture is composed of an encoder and decoder module.

Encoder: The encoder input consists of static components $C_{s,n} = \{w_n \cup \tau_n\}$ (location and expected serving time) and dynamic components $C_{d,n} = \{B_{\text{UAV}} \cup \mathcal{E}_n^{\text{ser}}\} \forall n = \{0 \dots 2N+1\}$ (UAV battery and energy requirements). These are processed through an embedding layer comprising a convolutional encoder, which maps the low-dimensional data to a high-dimensional space and produces embedded outputs C_{emb,s_n} and C_{emb,d_n} , with total inputs and outputs denoted as C_{in} and C_{emb} .

Decoder: For a policy ϕ , the probability that the UAV follows a sequence S conditioned on C_{emb} is defined as $P_{\phi}(S|C_{\text{emb}}) = \prod_{l=1}^L P(s_{l+1}|S_l, C_{\text{emb}})$, where S_l denotes the sequence up to l steps. At each step $l \in \{1, \dots, L\}$, the decoder generates the conditional probability distribution

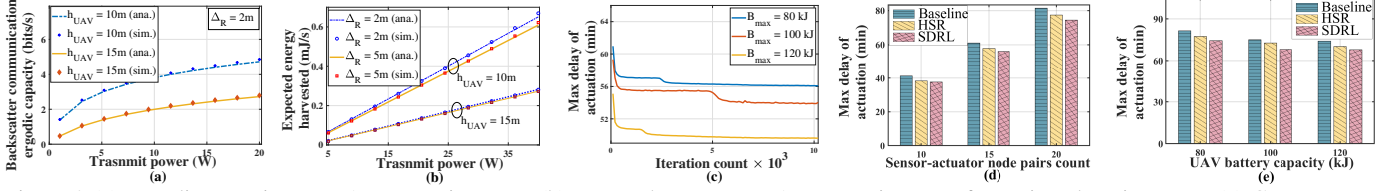


Figure 2 (a) Ergodic capacity vs. UAV transmit power. (b) Energy harvest vs. UAV transmit power for various location error. (c) Convergence of proposed SDRL algorithm; $N=15$. (d) Actuation delay, with battery size 80kJ. (e) Actuation delay for 20 sensor-actuator node pairs.

Table I Compute overhead for different sensor-actuator pair counts

| Method | Computation Time (ms) | | |
|----------|-----------------------|--------|--------|
| | N = 10 | N = 15 | N = 20 |
| Baseline | 0.1675 | 0.19 | 0.26 |
| HSR | 1504 | 2100 | 3080 |
| SDRL | 86.76 | 104.96 | 134.65 |

$P(s_{l+1}|S_l, C_{emb})$, that determines the agent action a_l . The optimal policy ϕ^* gives the optimal sequence S^* with probability 1. We aim to minimize the optimality gap between ϕ and ϕ^* .

In the l -th decoding step, the static information of the last visited index is encoded as $C_{emb,s_{l-1}}$, and Gated Recurrent Unit (GRU) processes this encoded data along with the hidden state h_{l-1} maintained by GRU, yielding $h_l, Y_l = \text{GRU}(h_{l-1}, C_{emb,s_{l-1}})$. The GRU output Y_l is then passed through the attention block which incorporates the importance of the locations at the current decoding step. This is achieved through the following operation: $a_l = \text{softmax}(v_{a1}^T \tanh(W_{a1} \cdot C_{emb}) + v_{a2}^T \tanh(W_{a2} \cdot Y_l))$. Using a_l , a context vector c_l , is generated as $c_l = \sum_{i=0}^{2N+1} (a_l^i C_{emb,i})$. The encoded output C_{emb} is then combined with c_l to produce κ_l as $\kappa_l = v_{c1}^T \tanh(W_{c1} \cdot C_{emb}) + v_{c2}^T \tanh(W_{c2} \cdot c_l)$. Before generating the action probability distribution, the masking vector M_l is applied to κ_l as $\kappa_l = \kappa_l + M_l$, where $M_{l,i} \in \{0, -\infty\} \forall i \in \{1, \dots, 2N+2\}$, enforces that the feasibility conditions in \mathcal{P}_1 , are satisfied. In addition, the actuator nodes are masked until all the sensor nodes are visited. Finally, the conditional probability distribution which determines the agent action a_l , is calculated as $P(s_{l+1}|S_l, C_{emb}) = \text{softmax}(\kappa_l)$. The process repeats until the termination condition is satisfied.

2) *Training*: The UAV policy network is trained using the REINFORCE algorithm to maximize expected cumulative reward, as outlined in Algorithm 1. The algorithm utilizes both an actor and a critic neural network to optimize the policy. The actor adjusts the policy parameters ϱ_a , increasing the likelihood of selecting actions that lead to positive outcomes. During each episode, the agent collects experience from the environment and computes cumulative reward R^b at the end, which serves as a performance estimate. Following the policy gradient approach, the policy parameters are then updated in the direction that increases the likelihood of actions that lead to higher rewards. The critic estimates the value function V^b and associated parameters ϱ_c are trained using stochastic gradient descent on the mean squared error (MSE) between predicted value and actual reward, helping to reduce variance in policy updates. To enhance learning efficiency, parallel learning and batch processing are used, addressing correlation between successive states and improving training stability.

Algorithm 1 Pseudo-code for SDRL training stage

```

1: Input: batch size  $\mathcal{B}$ , training dataset  $\mathcal{D}$ .
2: Output: Policy parameters  $\varrho_a, \varrho_c$ 
3: for itr = 1, 2, ..., i_tot do
4:   reset gradient  $d\varrho_a \leftarrow 0, d\varrho_c \leftarrow 0$ 
5:   Obtain training batch  $\{\mathcal{D}_{(itr \times \mathcal{B})+1}, \dots, \mathcal{D}_{(itr+1) \times \mathcal{B}}\}$ 
6:   for  $b = 1, 2, \dots, \mathcal{B}$  do
7:      $l \leftarrow 0$ 
8:     do
9:       compute  $P(s_{l+1}^b | S_l^b, C_{emb}^b)$ 
10:      choose action accordingly
11:       $l \leftarrow l + 1$ 
12:    while terminal state is reached
13:    compute reward  $R^b$ 
14:  end for
15:  calculate gradient
16:   $d\varrho_a \leftarrow \frac{1}{\mathcal{B}} \sum_{b=1}^{\mathcal{B}} (R^b - V^b) \nabla_{\varrho_a} \log(P^b(S|C_{emb}))$ 
17:   $d\varrho_c \leftarrow \frac{1}{\mathcal{B}} \sum_{b=1}^{\mathcal{B}} \nabla_{\varrho_c} (R^b - V^b)^2$ 
18:  Update  $\varrho_a$  and  $\varrho_c$ 
19: end for

```

VI. NUMERICAL RESULTS

The SDRL simulations are performed in Python with the PyTorch library on a system having an Intel Xeon W-2145 3.70 GHz processor, 64 GB RAM, and Nvidia Quadro P5000 GPU. Simulation parameters include: $m_1, m_2 = 2$; $\alpha_1, \alpha_2 = 1$; $D_{req} = 24-40$ Bytes; $E_{req} = 31-47$ mJ; $G_1 = 10$ dBi; $G_2 = 0$ dBi; $\tau_{BSS} = 180$ s; $h_{UAV} = 10$ m; $\lambda = 915$ MHz; $R_{max} = 150$ m; $\Delta R = 2$ m; $\alpha_P = 2$; $\varrho_1 = 0.000064$; $\varrho_2 = 0.02$; $P_c = 10.6$ μ W; $P_{EH,max} = 0.00492$ μ W; $\sigma_1^2, \sigma_2^2 = 10^{-9}$; constants $a = 0.24714$; $b = -1.5817 \times 10^{-5}$; $P_{tx}^{(a)} = 40$ W; $P_{tx}^{(s)} = 10$ W; $V_{UAV} = 9.8$ m/s; training size = 1.28×10^6 ; validation size = 10^3 ; convolution layers = 1; hidden size $h_s = 128$; learning rate = 5×10^{-4} ; dropout = 0.1; batch size = 256. Fig. 2(a) highlights the BSC ergodic capacity. The plot illustrates a positive correlation between ergodic capacity and transmit power, while elevated hovering altitudes result in a noticeable reduction in capacity due to increased path loss. The low data rate is primarily due to the multiplicative fading and dual pathloss in BSC link and compounded by UAV position imprecision. Fig. 2(b) illustrates the performance of the UAV energy transfer capability in terms of expected energy harvesting rate. The plot captures the effect of UAV altitude variation and localization inaccuracy. At lower elevations, the impact of localization uncertainty on performance is more pronounced. The results reflect the characteristics of monostatic BSC systems, where ultra-low-power operation is prioritized over high data rates. The observed performance aligns with low-rate, energy-constrained application scenarios, where small payloads and efficient energy usage

are key. The analytical expressions align well with Monte Carlo simulations, confirming both their accuracy and practical relevance within the operational context considered.

Fig. 2(c) shows the learning curves of the proposed algorithm for different battery capacities. The curves demonstrate eventual convergence to the long-term return. The maximum delay is observed with the lowest battery capacity. Figures 2(d) and 2(e) compare the proposed algorithm, the baseline solution, and the hybrid swap-and-reverse (HSR) algorithm. The baseline solution satisfies the feasibility constraints but does not necessarily yield an optimal result. In contrast, the HSR algorithm employs a 2-Opt-based approach [16], starting from a feasible solution and iteratively refining it by randomly selecting two indices to either swap nodes or reverse the subsequence between them. The update is retained if it improves the reward and maintains feasibility, progressively refining the solution. Fig. 2(d) demonstrates that the maximum actuation delay in the proposed approach is consistently lower than the benchmark schemes. Also, Fig. 2(e) shows a clear reduction in maximum actuation delay. For example, in a network with 20 sensor-actuator node pairs and 80 kJ UAV battery capacity, the proposed strategy reduces actuation delay by up to 182 s compared to benchmark methods across various scenarios.

The computational complexity of the proposed SDRL is $O(LN h_s^2 C_1)$, while that of HSR is $O(I_{max} N^2 C_2)$, where C_1, C_2 are constants, h_s is the hidden size, and I_{max} is the number of iterations. The computation times are compared in Table I, showing up to 22 times reduced overhead in the proposed approach compared to HSR, highlighting its efficiency. The proposed algorithm is robust to changes in node count or configuration, requiring no retraining and thereby ensuring effective adaptation to network variations.

VII. CONCLUSION

This paper introduced UAV and BSS-assisted WSN for minimizing actuation delay with imperfect UAV localization and wireless fading channel constraints. We derived closed-form expressions of BSC ergodic capacity for data collection and expected energy harvesting rate. Based on node requirements and the statistical properties of communication and energy harvesting, the node visit sequence of UAV was optimized using sequential DRL. Our results demonstrate that the proposed DRL-aided strategy offers consistently reduced actuation delay with multifold reduction in computation overhead compared to the benchmark approaches.

APPENDIX

A. Proof of (4): The pdf of the product of two independent not identical (i.n.i.d.) gamma random variables h_1 and h_2 is

$$f_{h_1 h_2}(z) = \frac{2}{\Gamma m_1 \Gamma m_2} \left(\frac{m_1 m_2}{\alpha_1 \alpha_2} \right)^{\frac{m_1 + m_2}{2}} z^{\frac{m_1 + m_2 - 2}{2}} K_{m_1 - m_2} \left(2 \sqrt{\frac{m_1 m_2 z}{\alpha_1 \alpha_2}} \right). \quad (.1)$$

The ergodic capacity is defined as

$$\frac{1}{\log 2} \int_{d_{\min}}^{d_{\max}} \int_0^{\infty} \log(1 + \gamma) f_{h_1 h_2}(x) f_d(d) dx dd \quad (.2)$$

where $\gamma = \eta d^{-2\alpha} x$, $d_{\min} = h_{UAV}$, $d_{\max} = \sqrt{\Delta_R^2 + h_{UAV}^2}$, using (.1) and [17, 07.34.03.0456.01]. Consider the following integral:

$$I_1 = C_1 \int_0^{\infty} G_{2,2}^{1,2} \left(\frac{11}{10} \middle| C_2 x \right) x^{\frac{m_1 + m_2 - 2}{2}} K_{m_1 - m_2} (2\sqrt{x}) dx. \quad (.3)$$

Using [18, Eq. (7.821.3)], and substituting in (.2), we have

$$\int_{d_{\min}}^{d_{\max}} C_3 G_{4,2}^{1,4} \left(\frac{1-m_1, 1-m_2, 1, 1}{1, 0} \middle| C_4 d^{-2\alpha} \right) \frac{d}{\sqrt{d^2 - h_{UAV}^2}} dd. \quad (.4)$$

Thereafter, using [18, Eq. (9.34.7)], and representing the output in terms of bivariate Fox-H function, we get (4).

B. Proof of (5): The expected energy harvesting rate is

$$\bar{P}_{EH} = \int_{\rho_1}^{\rho_2} (aP + b) f_{\bar{P}_{rx}}(P) dP + P_{EH, \max} \int_{\rho_2}^{\infty} f_{\bar{P}_{rx}}(P) dP \quad (.1)$$

where $\bar{P}_{rx} = E_d[P_{rx}] = G_1 (\lambda/4\pi)^2 E[d^{-\alpha P}] h_1 P_{tx}$. For $\alpha_P = 2$, after simple mathematical calculations, we get

$$E[d^{-\alpha P}] = \int_{d_{\min}}^{d_{\max}} d^{-\alpha P} f_d(d) dd = \frac{2}{h_{UAV} \Delta_R} \tan^{-1} \left(\frac{\Delta_R}{h_{UAV}} \right). \quad (.2)$$

Finally, substituting (.2) into (.1) and using [18, Eq. (8.350.1)], the closed-form expression (5) is obtained.

REFERENCES

- [1] N. Primeau *et al.*, "A review of computational intelligence techniques in wireless sensor and actuator networks," *IEEE Commun. Surveys Tuts.*, vol. 20, no. 4, pp. 2822–2854, 2018.
- [2] D. Mishra *et al.*, "Smart RF energy harvesting communications: challenges and opportunities," *IEEE Commun. Mag.*, vol. 53, no. 4, pp. 70–78, 2015.
- [3] T. Jiang *et al.*, "Backscatter communication meets practical battery-free internet of things: A survey and outlook," *IEEE Commun. Surveys Tuts.*, vol. 25, no. 3, pp. 2021–2051, 2023.
- [4] H. Fu *et al.*, "Analysis and optimization of age of information for area sensing," *IEEE Commun. Lett.*, vol. 28, no. 1, pp. 103–107, 2024.
- [5] K. Fizza *et al.*, "Age of data aware internet of things applications," in *IEEE CCNC*, 2022, pp. 399–404.
- [6] A. Nikkhah *et al.*, "Age of actuation in a wireless power transfer system," in *IEEE INFOCOM Wksp.*, 2023, pp. 1–6.
- [7] B. Chang *et al.*, "Age of information for actuation update in real-time wireless control systems," in *IEEE INFOCOM Wksp.*, 2020, pp. 26–30.
- [8] M. Hoang *et al.*, "Design of autonomous battery swapping for UAVs," in *IEEE/ASME Int. Conf. AIM*, 2024, pp. 353–358.
- [9] T. Cokyasari *et al.*, "Designing a drone delivery network with automated battery swapping machines," *Comput. Oper. Res.*, vol. 129, p. 105177, 2021.
- [10] Y. Li *et al.*, "Energy-efficient UAV-driven multi-access edge computing: A distributed many-agent perspective," *IEEE Trans. Commun.*, pp. 1–1, 2025.
- [11] Y. Zeng, J. Xu, and R. Zhang, "Energy minimization for wireless communication with rotary-wing UAV," *IEEE Trans. Wireless Commun.*, vol. 18, no. 4, pp. 2329–2345, 2019.
- [12] Y. H. Al-Badarneh *et al.*, "Performance analysis of monostatic multi-tag backscatter systems with general order tag selection," *IEEE Wireless Commun. Lett.*, vol. 9, no. 8, pp. 1201–1205, 2020.
- [13] D. Mishra, S. De, and D. Krishnaswamy, "Dilemma at RF energy harvesting relay: Downlink energy relaying or uplink information transfer?" *IEEE Trans. Wireless Commun.*, vol. 16, no. 8, pp. 4939–4955, 2017.
- [14] N. Mazyavkina *et al.*, "Reinforcement learning for combinatorial optimization: A survey," *Comput. Oper. Res.*, vol. 134, p. 105400, 2021.
- [15] M. Nazari *et al.*, "Reinforcement learning for solving the vehicle routing problem," in *Adv. Neural Inf. Process. Syst.*, vol. 31, 2018.
- [16] Y. Ren and V. Friderikos, "Path planning optimization based interference awareness for mobile robots in mmwave multi cell networks," *IEEE Trans. Veh. Technol.*, vol. 73, no. 9, pp. 13 639–13 650, 2024.

- [17] I. Wolfram, “Wolfram, research, mathematica edition: Version 10.0. champaign,” *Wolfram Research, Inc.*, 2010.
- [18] I. S. Gradshteyn and I. M. Ryzhik, “Table of Integrals, Series and Products,” *New York, NY, USA: Academic*, 2000.