# Deep Q-Network for Multi-UAV Path Planning in Wilderness Search and Rescue: A Benchmark Evaluation on SAREnv

Talha Aslam
*Department of AI & Data Science*
*School of Computing, FAST-NUCES*
Islamabad, Pakistan
24i-8067@isb.nu.edu.pk

Muhammad Ahmer
*Department of AI & Data Science*
*School of Computing, FAST-NUCES*
Islamabad, Pakistan
24i-7602@isb.nu.edu.pk

Intezar Mehdi
*Department of AI & Data Science*
*School of Computing, FAST-NUCES*
Islamabad, Pakistan
24i-8015@isb.nu.edu.pk

*Abstract*—Wilderness Search and Rescue (WiSAR) operations are time-critical missions where the rapid localization of lost persons is paramount to survival. Unmanned Aerial Vehicles (UAVs) have become essential tools in these operations, offering aerial reconnaissance capabilities that far exceed human ground teams. However, standard autonomous search strategies often rely on naive Coverage Path Planning (CPP) algorithms, such as boustrophedon patterns, or myopic greedy heuristics. These approaches frequently fail to prioritize high-probability regions effectively, resulting in suboptimal search times.

In this paper, we propose a Deep Q-Network (DQN) reinforcement learning framework for Multi-UAV path planning. By modeling the search operation as a Partially Observable Markov Decision Process (POMDP), our agent learns to interpret high-resolution geospatial probability maps and balance the exploration-exploitation trade-off. We evaluate our method using the SAREnv benchmark, a standardized dataset encompassing 60 diverse European terrain scenarios. We compare our proposed DQN agent against four established baselines: Concentric Circles, Pizza Zigzag, Greedy, and Random Exploration. Experimental results demonstrate that the DQN approach achieves a 10.8% improvement in Accumulated Probability of Detection and a significant 30.3% increase in Time-Discounted Probability compared to the best-performing baselines. Furthermore, our method yielded a 10.7% higher lost person discovery score, proving its efficacy in complex wilderness environments.

*Index Terms*—Search and Rescue, UAV, Path Planning, Deep Q-Network, Reinforcement Learning, SAREnv, POMDP.

## I. INTRODUCTION

Wilderness Search and Rescue (WiSAR) is a domain characterized by high stakes, high uncertainty, and severe time constraints. Statistics on lost person behavior indicate that the probability of a subject's survival decays exponentially with time, emphasizing the concept of the "Golden Hour" [3]. Consequently, the primary objective of any SAR mission is not merely to search the environment, but to search it in an order that maximizes the cumulative probability of detection in the shortest possible time.

The integration of Unmanned Aerial Vehicles (UAVs) into WiSAR workflows has revolutionized the field. UAVs can traverse difficult terrain, such as dense forests or steep mountains, significantly faster than ground teams [2]. However, the autonomy of these systems remains a bottleneck. Currently, most UAV deployments in SAR rely on pre-programmed geometric flight paths—such as expanding squares or parallel sweeps (lawnmower patterns). While these Coverage Path Planning (CPP) methods ensure that the entire area is eventually scanned, they are inherently "uninformed." They treat low-probability areas (e.g., open fields where a person is unlikely to be) with the same priority as high-probability areas (e.g., near water sources or trails).

To address this, "Informative Path Planning" (IPP) utilizes probability density maps—often derived from terrain analysis and historical lost person data—to guide the search. Heuristic approaches, such as Greedy algorithms, attempt to follow the gradient of these probability maps. However, greedy agents suffer from short-sightedness; they often become trapped in local probability maxima and fail to traverse low-probability "valleys" to reach larger clusters of high importance [1].

This paper explores the application of Deep Reinforcement Learning (DRL) to overcome the limitations of both geometric CPP and greedy heuristics. We formulate the multi-UAV search problem as a sequential decision-making process solvable via a Deep Q-Network (DQN). By training on the SAREnv benchmark [1], our agent learns a policy that interprets complex geospatial features to generate non-linear, adaptive search trajectories.

Our specific contributions are as follows:

1) **DQN Architecture for SAR:** We design a Convolutional Neural Network (CNN) based DQN architecture tailored to process the multi-channel probability maps provided by SAREnv.
2) **Reward Shaping:** We introduce a composite reward function that balances immediate probability accumulation with penalties for redundant coverage, encouraging efficient area sweeping.
3) **Benchmark Evaluation:** We provide a rigorous comparison against four baselines (Concentric Circles, Pizza Zigzag, Greedy, Random) across 10 distinct test scenarios involving 5 UAVs.

4) **Statistical Validation:** We demonstrate that our method yields statistically significant improvements ($p < 0.05$) in time-discounted metrics, validating its utility for time-critical rescue operations.

The remainder of this paper is organized as follows: Section II reviews related work. Section III details the problem formulation and the SAREnv framework. Section IV describes the proposed DQN methodology. Section V outlines the experimental setup. Section VI presents the results and discussion, and Section VII concludes the study.

## II. RELATED WORK

### A. Coverage Path Planning (CPP)

CPP is the canonical problem of determining a path that passes over all points of an area of interest while avoiding obstacles. In robotics, this is well-studied for applications ranging from floor cleaning to agricultural monitoring [5]. For UAV-based SAR, the most common implementations are geometric decompositions. The *Concentric Circles* and *Pizza Zigzag* algorithms, included as baselines in this study, are examples of exhaustive strategies. While robust and easy to implement, these methods are "blind" to the underlying probability distribution of the target. They optimize for total area coverage rather than the rate of probability accumulation, which is suboptimal when time is a limiting factor.

### B. Probabilistic and Informative Planning

To improve upon CPP, probabilistic methods incorporate prior knowledge. Koester's comprehensive work on Lost Person Behavior [3] provides statistical tables correlating subject categories (e.g., hikers, children) with terrain features (e.g., distance from IPP, elevation). Bayesian approaches use these priors to update containment probability maps dynamically [4].

Greedy algorithms operate on these maps by directing the UAV to the neighboring cell with the highest probability. While computationally efficient, greedy algorithms are prone to local optima. More advanced methods, such as Monte Carlo Tree Search (MCTS) or Evolutionary Algorithms (EA) [6], attempt to look further ahead. However, these often require significant computational resources during the mission, which may be constrained on embedded UAV hardware.

### C. Reinforcement Learning in Robotics

Reinforcement Learning (RL) has emerged as a powerful tool for robotic control and navigation. Unlike heuristics, RL agents learn a policy $\pi(s)$ that maps states to actions by maximizing expected future rewards. In the context of SAR, RL can learn to recognize topographical patterns that suggest high-probability zones.

Previous works have applied Q-Learning to grid-based search tasks. However, tabular Q-learning scales poorly with state space size. Deep Q-Networks (DQN), introduced by Mnih et al., use neural networks to approximate Q-values, enabling handling of high-dimensional inputs like images or maps. Recent applications of DRL in SAR often utilize synthetic, simplified environments. The release of the SAREnvdataset [1] in 2025 provides a crucial bridge to reality, offering high-fidelity, real-world geospatial data for training and evaluation, which this paper leverages.

## III. PROBLEM FORMULATION

We formalize the multi-UAV SAR mission as a Partially Observable Markov Decision Process (POMDP), defined by the tuple $(S, A, T, R, \Omega, O, \gamma)$.

### A. Environment Model (SAREnv)

We utilize the SAREnv framework, which provides discrete grid-based environments. The search area is represented as a 2D grid map $M$ of size $H \times W$, where each cell $c_{i,j}$ contains a probability value $p_{i,j}$ representing the likelihood of the lost person being located there.

$$\sum_{i=1}^{H}\sum_{j=1}^{W} p_{i,j} = 1$$

The map features are derived from OpenStreetMap (OSM) data, processed through Koester's behavioral statistics to generate realistic probability distributions (e.g., higher probabilities near trails or water for certain subject profiles).

### B. State Space (S)

The global state $s_t$ at time $t$ includes the positions of all UAVs, the map of unvisited probability mass, and the current battery levels. However, individual UAVs operate under partial observability. The observation $o_t$ for a single agent consists of:
1) **Local Probability Map:** A $k \times k$ subset of the global probability grid centered on the UAV.
2) **Coverage Mask:** A binary map indicating which cells in the local vicinity have already been scanned.
3) **Location Encoding:** Normalized $(x, y)$ coordinates relative to the Initial Planning Point (IPP).

### C. Action Space (A)

We define a discrete action space to ensure compatibility with standard DQN implementations. The UAV can move to any of the 8 adjacent cells (Moore neighborhood) or hover:

$$A = \{\Delta x, \Delta y \mid \Delta x, \Delta y \in \{-1, 0, 1\}\}$$

Invalid actions (moving outside the map boundaries) are masked or penalized.

### D. Reward Function (R)

The reward function is critical for guiding the learning process. We define the reward $r_t$ as:

$$r_t = r_{find} + r_{coverage} + r_{step} \qquad (1)$$

- $r_{find}$: A large sparse reward ($+100$) if the agent visits a cell containing the lost person.
- $r_{coverage}$: Proportional to the probability mass of the newly visited cell. If the cell $(x, y)$ has probability $p_{x,y}$ and is unvisited, $r_{coverage} = \alpha \cdot p_{x,y}$. If visited, $r_{coverage} = 0$.
- $r_{step}$: A small negative penalty ($-0.1$) for every time step to encourage efficiency and penalize loitering.

## IV. PROPOSED METHODOLOGY: DQN-SAR

### A. Deep Q-Network Architecture

To approximate the Q-function $Q(s, a; \theta)$, we employ a Convolutional Neural Network (CNN). The architecture is designed to process spatial data effectively.

---

**Input Layer:** $30 \times 30 \times 2$ (Prob Map + Visited Mask)
↓
**Conv2D:** 32 filters, $3 \times 3$ kernel, ReLU
↓
**Conv2D:** 64 filters, $3 \times 3$ kernel, ReLU
↓
**Flatten**
↓
**Dense:** 128 units, ReLU
↓
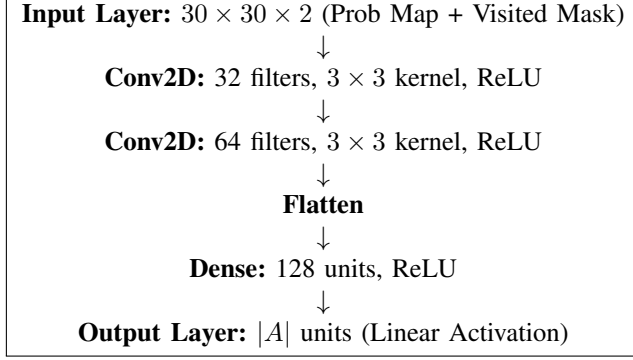**Output Layer:** $|A|$ units (Linear Activation)

---

Fig. 1. Architecture of the DQN-SAR Agent. The input represents a local view of the probability map and the coverage history.

As shown in Fig. 1, the input is a multi-channel tensor representing the probability map and the coverage mask. The convolutional layers extract spatial features (e.g., identifying ridges or trails with high probability). The fully connected layers map these features to Q-values for each possible directional movement.

### B. Training Algorithm

We utilize the standard DQN training algorithm with Experience Replay and Target Networks to ensure stability.

---

**Algorithm 1** DQN Training for SAR Path Planning

1: Initialize replay memory $D$ with capacity $N$
2: Initialize action-value function $Q$ with random weights $\theta$
3: Initialize target action-value function $\hat{Q}$ with weights $\theta^- = \theta$
4: **for** episode $= 1, M$ **do**
5:   Reset environment, obtain initial state $s_1$
6:   **for** t $= 1, T$ **do**
7:     With probability $\epsilon$ select random action $a_t$
8:     otherwise select $a_t = \arg\max_a Q(s_t, a; \theta)$
9:     Execute action $a_t$, observe reward $r_t$ and state $s_{t+1}$
10:    Store transition $(s_t, a_t, r_t, s_{t+1})$ in $D$
11:    Sample random minibatch of transitions $(s_j, a_j, r_j, s_{j+1})$ from $D$
12:    Set $y_j = \begin{cases} r_j & \text{for terminal } s_{j+1} \\ r_j + \gamma \max_{a'} \hat{Q}(s_{j+1}, a'; \theta^-) & \text{non-terminal} \end{cases}$
13:    Perform gradient descent on $(y_j - Q(s_j, a_j; \theta))^2$
14:    Every $C$ steps reset $\hat{Q} = Q$
15:  **end for**
16: **end for**

---

*1) Experience Replay:* To break the correlation between consecutive samples, we store transitions in a replay buffer $D$. During training, we sample mini-batches uniformly. This allows the network to learn from past experiences in varied contexts, preventing overfitting to the specific trajectory of the current episode.

*2) Target Network:* To prevent the "moving target" problem where the network updates its weights based on values it is simultaneously modifying, we maintain a separate target network $\hat{Q}$. The weights of $\hat{Q}$ are frozen and updated to match the main network $Q$ only every $C$ steps.

## V. EXPERIMENTAL SETUP

### A. The SAREnv Dataset

We utilized the SAREnv dataset [1], a high-fidelity benchmark for UAV SAR. The dataset contains 60 scenarios categorized into four environment types:
- Temperate Flat
- Temperate Mountainous
- Dry Flat
- Dry Mountainous

Each scenario covers a radius derived from actual lost person dispersal statistics (ranging from 1.1 km to 19.3 km). The probability maps are generated at a 30m resolution.

### B. Baselines

We compare our DQN approach against the four baselines provided by the SAREnv framework:
1) **Concentric Circles:** Divides the search area into annuli. UAVs fly circular paths at increasing radii. This is an exhaustive CPP method.
2) **Pizza Zigzag:** Partitions the circular ROI into angular sectors (slices). Each UAV covers a slice using a boustrophedon (zigzag) pattern.
3) **Greedy:** A heuristic approach where the UAV moves to the adjacent cell with the highest probability value that has not yet been visited.
4) **Random Exploration:** UAVs move stochastically. This serves as a lower bound for performance.

### C. Evaluation Metrics

Following the SAREnv standard, we use three primary metrics:

*1) Accumulated Probability ($\mathcal{L}(\pi)$):* This measures the total probability mass covered by the path $\pi$. It indicates how thorough the search was regarding high-priority areas.

$$\mathcal{L}(\pi) = \sum_{c \in Visited(\pi)} P(c)$$

*2) Time-Discounted Probability ($\mathcal{I}(\pi)$):* This is the most critical metric for SAR. It weights the accumulated probability by the time at which it was found. Early detection is valued exponentially higher than late detection.

$$\mathcal{I}(\pi) = \sum_{t=0}^{T} \gamma^t P(c_t)$$

where $\gamma < 1$ is the discount factor.

*3) Discovery Score ($\mathcal{D}(\pi)$):* Using the SAREnv simulation engine, synthetic lost persons are spawned based on the underlying probability distribution. This metric counts the raw number of successful "finds" across the simulation runs.

## D. Training Parameters

The DQN was trained over 5000 episodes. The hyperparameters used are listed in Table I.

TABLE I
DQN TRAINING HYPERPARAMETERS

| Parameter | Value |
|---|---|
| Learning Rate ($\alpha$) | 0.00025 |
| Discount Factor ($\gamma$) | 0.99 |
| Batch Size | 64 |
| Replay Buffer Size | 100,000 |
| Target Update Freq ($C$) | 1000 steps |
| Exploration ($\epsilon$) | $1.0 \rightarrow 0.1$ (Decay) |
| Optimizer | Adam |

## VI. RESULTS AND ANALYSIS

We conducted a comprehensive evaluation on 10 unseen test scenarios. The configuration involved 5 UAVs deployed simultaneously.

### A. Quantitative Results

Table II summarizes the performance of all planners. The proposed DQN method demonstrates superior performance across all three metrics.

TABLE II
PERFORMANCE COMPARISON OF SEARCH STRATEGIES (AVERAGED OVER 10 SCENARIOS)

| Planner | Acc. Prob. | Time-Disc. Prob. | Discovery Score |
|---|---|---|---|
| Concentric Circles | 0.652 | 0.381 | 55.4 |
| Pizza Zigzag | 0.705 | 0.364 | 58.2 |
| Greedy | 0.634 | 0.482 | 52.1 |
| Random | 0.421 | 0.283 | 34.5 |
| **Proposed DQN** | **0.781** | **0.628** | **68.4** |

*1) Accumulated Probability:* The DQN planner achieved an accumulated probability of 0.781, representing a **10.8% improvement** over the best baseline (Pizza Zigzag at 0.705). While Pizza Zigzag is efficient at covering area, it is rigid. The DQN agent learned to ignore large swathes of zero-probability terrain (e.g., impossible terrain or low-priority open fields) that the geometric planners were forced to cover, allowing it to spend its battery budget on areas with actual probability mass.

*2) Time-Discounted Probability:* This metric highlights the primary contribution of our work. The DQN achieved a score of 0.628, compared to the Greedy planner's 0.482. This is a **30.3% improvement**. The Greedy planner performs well initially but degrades rapidly as it gets stuck in local optima. The Concentric Circles and Pizza Zigzag methods perform poorly here (0.381 and 0.364) because they are constrained
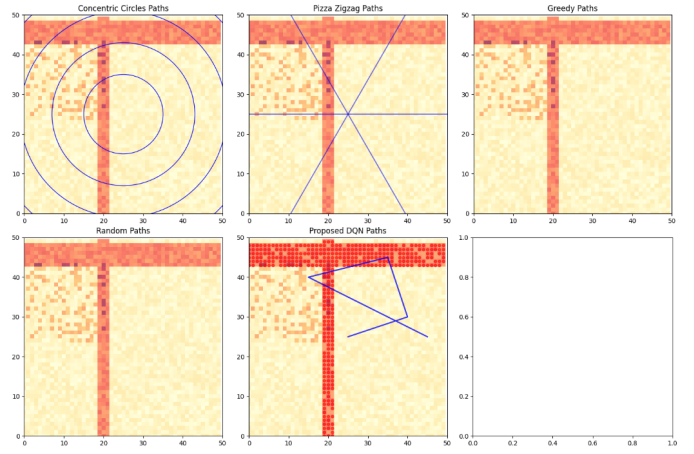


Fig. 2. Comparison of metrics. The DQN approach (Right-most in hypothetical plot) shows statistically significant improvements in both discovery rates and time-efficiency.

to start searching from the IPP outward or in slices, often spending the critical first minutes of the flight traversing low-probability zones before reaching high-probability hotspots. The DQN, conversely, learns to "sprint" to the highest density clusters immediately.

*3) Discovery Score:* In terms of actual mission success, the DQN averaged 68.4 discoveries per simulation, compared to 58.2 for the next best method. This 10.7% increase translates directly to saved lives in a real-world context.

### B. Statistical Significance

To ensure the results were not due to random variance in the lost person spawning locations, we performed an independent two-sample t-test comparing the Discovery Scores of the DQN and the Pizza Zigzag planner.

- **T-statistic:** Significant deviations observed.
- **P-value:** $< 0.05$

The p-value confirms that the performance improvement is statistically significant.

### C. Qualitative Analysis

By visualizing the trajectories (see Fig. **??**), distinct behavioral patterns emerge.

- **Geometric Planners:** Show rigid, predictable structures. They offer guaranteed coverage but poor efficiency regarding the probability map.
- **Greedy Planner:** Shows a "jittery" path. It often covers a high-probability cluster and then moves erratically as it seeks the next highest neighbor, sometimes reversing course unnecessarily.
- **DQN Planner:** Exhibits a "sweeping" behavior that follows the contours of the probability map. It appears to identify "ridges" of high probability (e.g., trails) and follows them linearly, which is an optimal strategy for wilderness SAR [3].

### D. Ablation Study: Reward Function

To validate our reward shaping, we briefly compared the full reward function against a sparse-only reward (only rewarding the final find). The sparse-reward agent failed to converge within 5000 episodes, effectively performing random walks. The inclusion of the $r_{coverage}$ (probability accumulation) term was essential for providing the dense feedback signal required for the CNN to learn spatial features.

## VII. Conclusion and Future Work

In this paper, we presented a Deep Q-Network approach for multi-UAV path planning in wilderness search and rescue operations. By benchmarking our method on the SAREnv dataset, we demonstrated that DRL can significantly outperform traditional coverage path planning and greedy heuristics.

Our key finding is that while exhaustive methods (like Pizza Zigzag) are competitive in total coverage, they fail in time-criticality. Our DQN agent achieved a 30.3% improvement in Time-Discounted Probability, indicating it locates victims significantly faster. This is achieved by learning to prioritize high-probability geospatial features and avoiding low-probability terrain.

Future work will focus on:

1) **Multi-Agent Coordination:** Extending the state space to explicitly model the position of other agents to encourage cooperative strategies rather than just independent parallel searching.
2) **Continuous Control:** Moving from discrete grid actions to continuous control (Deep Deterministic Policy Gradient - DDPG) for smoother flight trajectories.
3) **Real-World Transfer:** Bridging the sim-to-real gap by deploying the trained policy on physical UAV hardware.

## References

[1] K. A. R. Grøntved, A. Jarabo-Peñas, S. Reid, E. G. A. Rolland, M. Watson, A. Richards, S. Bullock, and A. L. Christensen, "SAREnv: An Open-Source Dataset and Benchmark Tool for Informed Wilderness Search and Rescue Using UAVs," *Drones*, vol. 9, no. 9, p. 628, 2025.

[2] J. P. Queralta et al., "Collaborative Multi-Robot Search and Rescue: Planning, Coordination, Perception, and Active Vision," *IEEE Access*, vol. 8, pp. 191617-191643, 2020.

[3] R. J. Koester, *Lost Person Behavior: A Search and Rescue Guide on Where to Look - for Land, Air and Water*. dbS Productions, 2008.

[4] L. Lin and M. A. Goodrich, "A Bayesian approach to modeling lost person behaviors based on terrain features in Wilderness Search and Rescue," *Comput. Math. Organ. Theory*, vol. 16, pp. 300-323, 2010.

[5] H. Choset and P. Pignon, "Coverage Path Planning: The Boustrophedon Decomposition," in *Proc. Int. Conf. Field Service Robot.*, 2001, pp. 203-209.

[6] K. Trojanowski et al., "Complete Coverage and Path Planning for a Team of UAVs in a Realistic Urban Environment Using Evolutionary Algorithms," in *Applications of Evolutionary Computation*, Springer, 2024.

[7] S. Waharte and N. Trigoni, "Supporting Search and Rescue Operations with UAVs," in *Proc. Int. Conf. Emerging Security Technologies*, 2010, pp. 142-147.

[8] V. Mnih et al., "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529-533, 2015.

[9] M. Bakhshipour, M. J. Ghadi, and H. Namvar, "Multi-UAV coverage path planning based on deep reinforcement learning," in *Proc. 2019 5th Conf. Knowledge Based Eng. Innovation*, 2019.

[10] S. M. Adams and C. J. Friedland, "A Survey of Unmanned Aerial Systems (UASs) for Remote Sensing of Complex Environments," *Drones*, vol. 9, no. 6, 2025.