

Multi-Agent Deep Reinforcement Learning For Persistent Monitoring With Sensing, Communication, and Localization Constraints

Manav Mishra, Prithvi Poddar, Rajat Agrawal, Jingxi Chen, Pratap Tokekar, and P.B. Sujit

Abstract—Determining multi-robot motion policies for persistently monitoring a region with limited sensing, communication, and localization constraints in non-GPS environments is a challenging problem. To take the localization constraints into account, in this paper, we consider a heterogeneous robotic system consisting of two types of agents: anchor agents with accurate localization capability and auxiliary agents with low localization accuracy. To localize itself, the auxiliary agents must be within the communication range of an anchor, directly or indirectly. The robotic team’s objective is to minimize environmental uncertainty through persistent monitoring. We propose a multi-agent deep reinforcement learning (MARL) based architecture with graph convolution called Graph Localized Proximal Policy Optimization (GALOPP), which incorporates the limited sensor field-of-view, communication, and localization constraints of the agents along with persistent monitoring objectives to determine motion policies for each agent. We evaluate the performance of GALOPP on open maps with obstacles having a different number of anchor and auxiliary agents. We further study (i) the effect of communication range, obstacle density, and sensing range on the performance and (ii) compare the performance of GALOPP with area partition, greedy search, random search, and random search with communication constraint strategies. For its generalization capability, we also evaluated GALOPP in two different environments – 2-room and 4-room. The results show that GALOPP learns the policies and monitors the area well. As a proof-of-concept, we perform hardware experiments to demonstrate the performance of GALOPP.

I. INTRODUCTION

Visibility-aware persistent monitoring (PM) problem involves continuous surveillance of a bounded environment by a single agent or a multi-agent system considering limited field-of-view (FOV) constraints into account [1]–[9]. Various applications, like search and rescue, border patrol, critical infrastructure, etc., require persistent monitoring for timely information. Ideally, persistent monitoring requires spatial and temporal separation of a team of robots in a larger environment to cooperatively carry out effective surveillance. The problem becomes complex as the multi-robot systems are subjected to limited sensing range, communication range, and localization constraints due to non-GPS environments. In this paper, we study the problem of determining motion planning policies for each agent in a multi-agent system for persistently monitoring a given environment, considering all

Manav Mishra, Prithvi Poddar, Rajat Agrawal, and P.B. Sujit are with the Department of Electrical Engineering and Computer Science, IISER Bhopal, Bhopal, India – 462038. Their work was supported by the Prime Minister Research Fellowship.

Jingxi Chen and Pratap Tokekar are with the Department of Computer Science, University of Maryland, College Park, United States – MD 20742. Their work was supported by NSF under grant No. 1943368, ONR under grant number N00014-18-1-2829, and an Amazon Research Award.

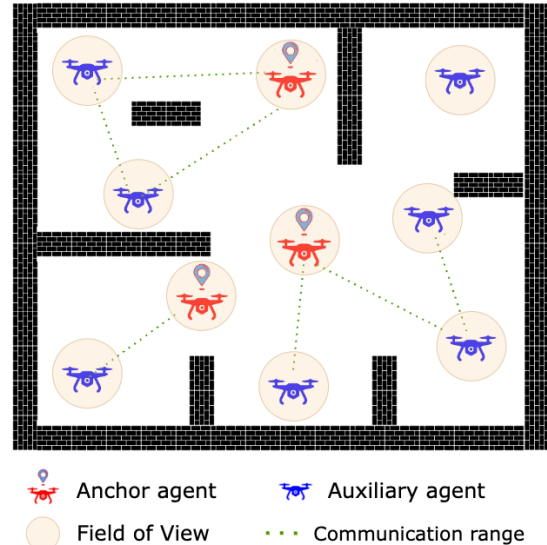


Fig. 1: Persistent monitoring in a 2-D environment using a team of anchor and auxiliary agents with FOV, localization, and communication range constraint.

the constraints using a graph communication-based multi-agent deep reinforcement learning (MARL) framework.

Generating motion policies for each agent using deterministic strategies becomes challenging due to the above constraints, as the agents require complete information about all possible interactions with information sharing among the agents. Hence, developing alternate strategies for multi-agent systems to learn to monitor complex environments is imperative.

We consider a scenario where a team of robots equipped with a limited field-of-view (FOV) sensor and limited communication range is deployed to persistently monitor a GPS-denied environment, as shown in Figure 1. As the environment does not support GPS, one can deploy agents that have expensive sensors such as tactical grade IMUs or cameras/LIDARs in conjunction with high computational power to carry out onboard SLAM for accurate localization with very low position uncertainty; these agents are called anchor agents. However, such a system becomes highly expensive for deployments. On the other hand, we can deploy agents with low-grade IMUs that are cheaper but exhibit high drift resulting in poor localization accuracy; these agents are called auxiliary agents. Auxiliary agents can be used in conjunction with external supporting localization units (like UWB ranging or cooperative localization [10]–[12])

to reduce localization uncertainty so that they are helpful in performing the coverage. Hence, as a trade-off, in this paper, we consider a robotic team consisting of anchor and auxiliary agents to monitor a region persistently.

The auxiliary agents can localize using the notion of cooperative localization by communicating with the anchor agents directly or indirectly through other auxiliary agents and hence have reduced uncertainty in their positional beliefs. As the auxiliary agents need to be in communication with anchor agents, their motion is restricted, which can result in lower monitoring performance as some areas may not be covered.

However, intermittent connection with the anchor agents will enable auxiliary agents to recover from the localization uncertainty while maintaining coverage across all regions [13]. This conflicting objective of monitoring the complete area while periodically maintaining connectivity from the anchor agents makes the problem of determining persistent monitoring strategies for the agents challenging.

Simultaneously addressing numerous constraints is pivotal in this context due to the intricacies involved in persistent monitoring. In GPS-denied scenarios, conventional methods relying on GPS or beacons fall short. Moreover, assuming complete communication between agents is often unrealistic, making it imperative to develop strategies that operate under partial or intermittent connectivity.

Consider applications such as search and rescue in urban/remote regions affected by flooding. In these scenarios, UAVs can be used to monitor and enhance situational awareness. However, the cameras have a finite field-of-view, and the communication range with the base station is also limited. In addition to these two constraints, acquiring GPS becomes challenging in cloudy weather conditions. Therefore, it is necessary to determine motion policies for the agents while considering all these constraints.

In this paper, we propose Graph Localized Proximal Policy Optimization (GALOPP), a multi-agent proximal policy optimization [14] algorithm coupled with a graph convolutional neural network [15] to perform persistent monitoring with such heterogeneous agents subject to sensing, communication, and localization constraints. The persistent monitoring environment is modeled as a two-dimensional discrete grid, and each cell in the grid is allocated a negative reward. When a cell is within the sensing range of any agent, then the reward value reduces to zero. Otherwise, the negative reward accumulates over time. Thus, the agents must learn their motion strategy to minimize the net reward accumulated over time, showing efficient persistent monitoring. We consider PPO in GALOPP because it is known for its stability, high sample efficiency, and resistance to hyperparameter tuning [14].

The approach presented in this paper addresses the challenge of balancing exploration for monitoring while adhering to localization constraints. Its novelty lies in the development of a model for guiding multiple decentralized agents in conducting monitoring in an uncertain environment. The primary contribution involves designing a decision-making strategy

based on Multi-Agent Reinforcement Learning for effective area surveillance while considering sensing, communication, and localization constraints.

The main contributions of this paper are:

- Development of a multi-agent deep reinforcement learning algorithm (GALOPP) for persistently monitoring a region considering the limited sensing range, communication, and localization constraints into account.
- Evaluating the performance of GALOPP for varying parameters – sensing area, communication ranges, the ratio of anchor to auxiliary agents, obstacle density, and centralized vs decentralized map sharing.
- Comparing the performance of GALOPP to baseline approaches, namely, random search, random search with ensured communication, greedy search, and area partitioning with lawn-mower sweeping strategy.

The rest of the paper is structured as follows. In Section II, we provide a review of the existing literature on this problem. In Section III, we define the persistent monitoring problem with multiple agents. In Section IV, we describe the GALOPP architecture and we evaluate the performance of GALOPP in Section V. In Section VI, the proof-of-concept of GALOPP performance using a team of nanocopters is described and we conclude in Section VII.

II. RELATED WORK

The persistent monitoring problem can be considered as a persistent area coverage or as a persistent routing problem visiting a set of targets periodically. Under persistent area coverage problem, one can consider the mobile variant of the Art Gallery Problem (AGP) [16] where the objective is to find the minimum number of guards to cover the area. There are several variants on AGP for moving agents under visibility constraints [17]. An alternative way for coverage is to use cellular decomposition methods, where the area can be decomposed into cells, and the agents can be assigned to these cells for coverage [18], [19]. In AGP and its variants, the visibility range is infinite but restricted by environmental constraints such as obstacles or boundaries.

In addressing the persistent routing problem, one can approach it using different variants of multiple Watchman Route Problem (n-WRP) [20]. In these approaches, the goal is to find a route for each agent for monitoring while minimizing the latency in visit time. Yu et al. [1] propose a method for monitoring events with stochastic arrivals at multiple stations by combining stochastic modeling and control theory to optimize real-time monitoring and resource utilization. Tokekar et al. [17] propose a novel method for persistent monitoring using robot teams that employs a coverage path planning algorithm accounting for visibility constraints to optimize coverage and path length trade-off. Wang et al. [8] propose a method for cooperative persistent surveillance on a road network using multiple Unmanned Ground Vehicles (UGVs) with detection capabilities. Lin et al. [4] apply a decentralized control algorithm to consider agents' dynamics and monitoring constraints in real-time, enabling the solution of the problem by finding the optimal

trajectory for each agent. Washington et al. [21] propose an RSVI algorithm for real-time drone surveillance policy optimization considering battery life and charging constraints, balancing the trade-off between surveillance coverage and energy consumption. Maini et al. [22] propose a coverage algorithm that considers the visibility constraints of robots to monitor linear features such as roads, pipelines, and power lines on terrains with obstacles. Mersheeva et al. [7] develop a framework for multi-UAV monitoring with priorities that address the efficient allocation of UAVs while considering resource limitations.

The above-cited approaches assume agents with either unrestricted sensing and/or communication and have full localization. In contrast, our context involves agents with limited access to information in these aspects. Devising a control policy that allows for persistent area coverage while respecting localization constraints is challenging. In the absence of GPS, the agent is unaware of its true position, and the estimated position covariance steadily increases over time. Therefore, the above approaches cannot be applied directly, and modifying them to accommodate localization constraints is difficult. The interconnection of exploration and localization objectives increases the problem's complexity.

Another approach is to learn from the environment to determine agent paths while considering the sensing, communication, and localization constraints. Reinforcement learning can be one such learning-based approach that can learn to determine paths for multiple agents while considering all the constraints. Multi-agent reinforcement learning (MARL) based path planning literature focuses on developing efficient and effective algorithms for multi-agent systems on cooperative multi-agent tasks covering a broad spectrum of applications [23]–[28]. Blumenkamp et al. [29] study inter-agent communication for self-interested agents in cooperative path planning but do not account for localization constraints and assume complete connectivity throughout. Omidshafiei et al. [23] formalize the concept of MARL under partial observability, which applies to scenarios with limited sensing range. Chen et al. [30] developed a method to find trajectories for agents to cover an area continuously but with the assumption that all agents have full access to the environment due to unrestricted communication access among agents. In the above articles, the problem of determining motion policies for the agents considering the localization, sensing, and communication range constraints jointly has not been adequately addressed. In this work, through GALOPP, we address the problem of persistent monitoring considering all three constraints using a deep reinforcement learning framework.

III. PROBLEM STATEMENT

A. Persistent monitoring problem

We consider the persistent monitoring problem in a 2D grid world environment $G \subseteq \mathbb{R}^2$ of size $A \times B$. Each grid cell $G_{\alpha\beta}$, $1 \leq \alpha \leq A$ and $1 \leq \beta \leq B$, has a reward $R_{\alpha\beta}(t)$ associated with it at time t . When the cell $G_{\alpha\beta}$ is within

the sensing range of an agent, then $R_{\alpha\beta}(t) \rightarrow 0$, otherwise, the reward decays linearly with a decay rate $\Delta_{\alpha\beta} > 0$. We consider negative reward as it refers to a penalty on the cell for not monitoring. At time $t = 0$, $R_{\alpha\beta}(t) = 0$, $\forall(\alpha, \beta)$ and

$$R_{\alpha\beta}(t+1) = \begin{cases} \max\{R_{\alpha\beta}(t) - \Delta_{\alpha\beta}, -R_{\max}\} & \text{if } G_{\alpha\beta} \text{ is not monitored at time } t \\ 0 & \text{if } G_{\alpha\beta} \text{ is monitored at time } t, \end{cases} \quad (1)$$

where R_{\max} refers to the maximum penalty a grid cell can accumulate so that the negative reward $R_{\alpha\beta}$ is bounded.

The objective of the persistent monitoring problem is to find a policy for the agents to minimize the neglected time, which in turn, maximizes the total accumulated reward by G over a finite time T . The optimal policy is given as

$$\pi^* = \arg \max_{\pi} \sum_{t=0}^T \left[\sum_{\alpha=1}^A \sum_{\beta=1}^B R_{\alpha\beta}^{\pi}(t) \right], \quad (2)$$

where π^* is an optimal global joint-policy that dictates the actions of the agent in a multi-agent system, and $R_{\alpha\beta}^{\pi}$ is the reward obtained by following a policy π .

Problem: Given a 2D grid world environment G , determine a joint-optimal policy π^* for N agents to minimize the neglect time at each cell in G taking sensing range, communication range and localization constraints into account.

B. Localization for Persistent Monitoring

The grid G consists of N -agents to perform the monitoring task. The agents have a communication range ρ . At every time step, a connectivity graph $\mathcal{G} = \langle \mathcal{V}, \mathcal{E} \rangle$ is generated between the agents. An edge connection e_{ij} is formed between agents i and j , if $dist(i, j) \leq \rho$, where $dist(i, j)$ is the Euclidean distance between agents i and j . The connectivity of any agent with an anchor agent is checked by using the Depth-First Search (DFS) algorithm. Each agent estimates its position using Kalman Filter (KF). The anchor agents have high-end localization capabilities; hence, the position uncertainty is negligible. However, the auxiliary agents can localize accurately if they are connected to an anchor agent, either directly or indirectly (multi-hop connection). As we do not consider any loss of information or cost associated with communication to an anchor agent, an auxiliary agent achieves localization upon observing an anchor agent. Consequently, a multi-hop communication to an anchor agent can localize another auxiliary agent not in direct communication contact with an anchor agent [12].

An agent located at position (α, β) has a field of view that covers a square region with dimensions $g \times g$, where $g = 2\ell + 1$, and the agent can sense ℓ cells in all cardinal directions.

As the anchor agents are accurately localized, they can update the rewards $R_{\alpha\beta}(t)$ in the grid world G , that is, set $R_{\alpha\beta}(t) = 0$. The auxiliary agents connected to the anchor agents either directly or via multi-hop connections can also update the rewards $R_{\alpha\beta}(t) = 0$. However, those auxiliary agents that are disconnected from the anchor agents can observe the world but cannot update the rewards due

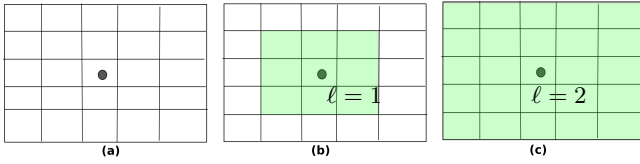


Fig. 2: Sensing range of the agent (a) agent position (b) When sensing range $\ell = 1$, the cells that the vehicle can sense $g = 3 \times 3$. (c) When $\ell = 2$, the sensing grid becomes $g = 5 \times 5$

to localization uncertainty associated with an increase in the covariance of the vehicle. When the vehicle reconnects with the anchor vehicle network, its uncertainty is reduced, and it can update the rewards. The world that the auxiliary agent observes during disconnection is not considered for simplicity.

An interesting aspect of solving Equation (2) to determine policies for the agents is that it does not explicitly assume that the graph network is always connected. Although a strict connectivity constraint increases the global positional belief of the entire team, it reduces the ability of the team to monitor any arbitrary region persistently due to the communication-constrained motion of the agent. Intermittent connectivity of agents leads to a better exploration of the area allowing more flexibility [31], [32]. The auxiliary agents, once disconnected, do not contribute to the net rewards obtained by the team. Since the objective is to find a policy that maximizes the rewards, the problem statement enables the agents to learn that connectivity increases the rewards, so they should be connected. The connectivity constraints are indirectly implied through rewards and not hard-coded into the agent decision-making policy. We abstract the localization constraints through the connectivity graph G during decision-making, which is detailed in section IV.

C. Using Kalman Filter for state estimation

In a cooperative localization (CL) setting, one way an auxiliary agent can localize is by observing an anchor agent. We assume that all the agents know their starting position accurately.

To handle the position uncertainties, we apply a Kalman Filter (KF) [33] to update its state mean and covariance. The KF propagates the uncertainty in the position of the auxiliary agent as long as it is unlocalized, and upon localization, the agent is made aware of its true location. The motion model of the auxiliary agent is,

$$\mu_{t+1} = A_t \mu_t + B_t u_t + \epsilon_t, \quad (3)$$

where μ_t and μ_{t+1} are the positions of the agents at time t and $t + 1$ respectively, ϵ_t is a random variable representing the error in the prediction, drawn from a normal distribution with zero mean and covariance R_t , $A_t = B_t = I_{2 \times 2}$, and u_t is the control input at time t . Upon observing another agent, the observation model can be formulated as,

$$z_t = C_t \mu_t + \delta_t, \quad (4)$$

$$C_t = \begin{bmatrix} \frac{a}{x_g - a} & 0 \\ 0 & \frac{b}{y_g - b} \end{bmatrix}. \quad (5)$$

Here $z_t = \begin{bmatrix} a \\ b \end{bmatrix}$ is the relative position of the observed agent in the context of the KF observation model, where $a = x_g - x$ and $b = y_g - y$. Here, (x_g, y_g) is the global position of the observed agent and (x, y) is the current position of the observing agent. δ_t is the error in the observation that is drawn from a normal distribution with 0 mean and covariance Q_t . Given the motion and observation models, we can write the KF algorithm as mentioned in Algorithm 1.

Algorithm 1 KF ($\mu_{t-1}, \Sigma_{t-1}, \mathbf{u}_t, z_t, \text{gotObservation}$)

$$\begin{aligned} \bar{\mu}_t &= A_t \mu_{t-1} + B_t \mathbf{u}_t \\ \bar{\Sigma}_t &= A_t \Sigma_{t-1} A_t^T + O_t \end{aligned}$$

if gotObservation \rightarrow **True** **then** $K_t = \bar{\Sigma}_t C_t^T (C_t \bar{\Sigma}_t C_t^T + Q_t)^{-1}$

$$\begin{aligned} \mu_t &= \bar{\mu}_t + K_t (z_t - C_t \bar{\mu}_t) \\ \Sigma_t &= (I - K_t C_t) \bar{\Sigma}_t \end{aligned}$$

return μ_t, Σ_t

else return $\bar{\mu}_t, \bar{\Sigma}_t$

Based on the environment model, vehicle motion, and localization model, we introduce our proposed GALOPP multi-agent reinforcement learning architecture in the next section.

IV. GRAPH LOCALIZED PPO - GALOPP

The multi-agent persistent monitoring task requires every individual agent to compute its policies using its own and the neighboring agents' observations. This makes computing policy for an agent a non-stationary problem that can be tackled using either a centralized or a decentralized algorithm. A centralized approach will comprise a single actor-critic network to determine the agents' policy. Such an algorithm is faster to train and execute but is not scalable to many agents. The decentralized approach overcomes these shortcomings by assigning individual actor networks to each agent. However, training multiple networks can be computationally expensive. In this paper, we utilize the Centralized Training and Decentralized Execution (CTDE) [34] strategy. This helps in retaining the computational efficiency of centralized actor-critic and the robustness of decentralized actors.

A. Architectural overview

The complete GALOPP pipeline with environmental interaction is shown in Figure 3, while the GALOPP architecture details are shown in Figure 4a. The GALOPP architecture consists of a multi-agent actor-critic model that implements Proximal Policy Optimization (PPO) [14] to determine individual agent policies. Multi-agent PPO is preferred over other policy gradient methods to avoid having large policy updates and achieve better learning stability in monitoring tasks. It also has better stability, high sample efficiency, and resistance to hyperparameter tuning [35].

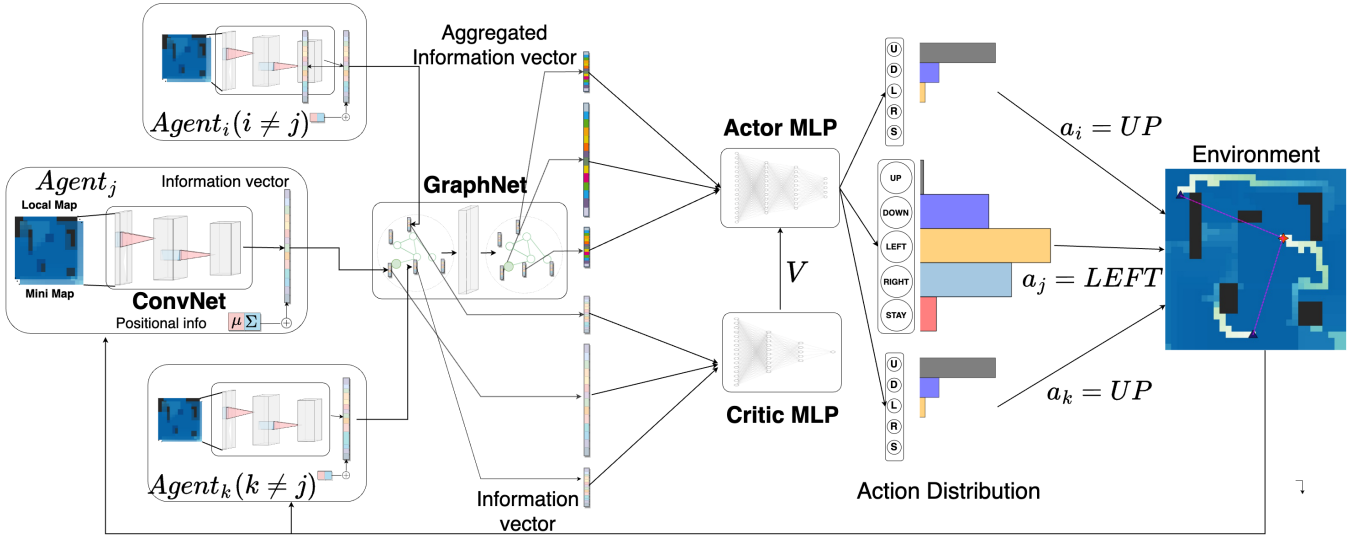


Fig. 3: Complete pipeline of the proposed GALOPP model with environmental interaction. The observations from each agent are processed by the ConvNet, and the generated embeddings are passed to the GraphNet following the communication graph formed among the agents. The GraphNet processes the input embeddings and generates aggregated information vectors that are passed through the actor network. The actor network generates a probability distribution over the possible actions for each agent, and the agents execute the actions having the highest probability. The critic provides feedback to the actor about the actions' expected value with respect to achieving the RL objective.

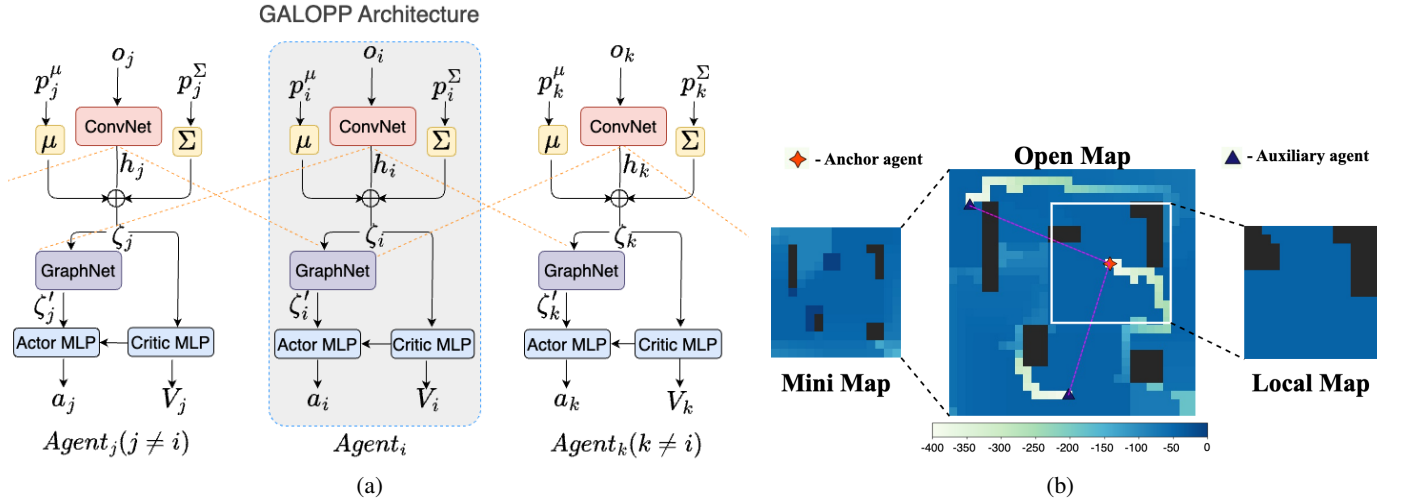


Fig. 4: (a) Schematic representation of GALOPP architecture. Each agent block of the architecture represents an actor-critic model. (b) The mini-map is the image of the environment G , resized to $g \times g$. The local map is a $g \times g$ slice of the environment G centered around the agent. The mini-map and local map are concatenated together to form the input o_i for agent i .

The observation of agent i is denoted as o_i , comprising a 2-channel image: the first channel, termed local map, represents the locally observed visibility map, while the second channel consists of an independently maintained global map version, compressed to match the dimensions of the local map (as depicted in Fig. 4b) and referred to as mini-map. The local map values depict a binary map indicating obstacle presence within the agent's geometric visibility constraint. The global map values represent the reward value heatmap of each cell in the grid, which is subsequently compressed to

form a mini-map. This image is passed through a Convolutional Neural Network (ConvNet) [36] to generate individual embeddings h_i for each agent, which is then augmented with agent i 's positional mean μ_i and covariance Σ_i , as shown in Figure 4a. This is the complete information ζ_i of the agent's current state. This information vector ζ_i forms the node embedding of the graph G . It is then processed by a Graph Convolutional Network (GraphNet) [15] that enforces the relay of messages in the generated connectivity graph G to ensure inter-agent communication. The decentralized actors

then use the embeddings generated by GraphNet to learn the policy, while a centralized critic updates the overall value function of the environment. The model is trained end-to-end for the persistent monitoring problem. The local computation involves - updating the local map, the mean and covariance of the position, and updating each agent’s maintained global map. The central computation is the computation of the joint policy for the multi-agent RL problem. The components of the GALOPP architecture are described in the below subsections.

B. Embedding extraction and message passing

The GALOPP model inputs the shared global reward values in the 2D grid. The observation of an agent i at time t is the set of cells that are within the sensing range (termed as the local map) and also a compressed image of the current grid (termed as mini-map) with the pixel values equal to the penalties accumulated by the grid cells [37]. Each agent has a separate copy of the mini-map. Each agent updates the copy of their mini-map, and the monitoring awareness is updated through inter-agent connectivity. Figure 6 illustrates a representation of how the decentralized map is updated. The connected agents compare and aggregate the global map at each time step for a network graph by taking the element-wise maximum for each grid cell $G_{\alpha\beta}$ in the environment. The element-wise maximum value of each grid cell is shared among the connected agents. The mini-map is resized to the shape of the local map of the agent and then concatenated to form a 2-channel image (shown in figure 4b). This forms the sensing observation input o_i for the model at time t . The ConvNet converts the observation o_i into a low-dimensional feature vector h_i termed the embedding vector. The positional mean μ_i and covariance matrix Σ_i of agent i is then flattened, and their elements are concatenated with h_i to generate a new information vector ζ_i (as shown in figure 4a).

The agents are heterogeneous agents (anchor and auxiliary) where the localization information is a parameter aggregated in the graph network component of GALOPP. An agent’s aggregated information vector depends on the current position in the environment, the generated message embedding, and the localization status of each neighboring agent.

GraphNet transfers the information vector ζ_i to all agents within the communication graph. The agents take in the weighted average of the embeddings of the neighborhood agents. The basic building block of a GraphNet is a graph convolutional layer, which is defined as [15]:

$$H^{(k+1)} = \sigma(A_g H^{(k)} W^{(k)}), \quad (6)$$

where $H^{(k)}$ is the feature matrix of the k -th layer, with each row representing a node in the graph and each column representing a feature of that node. A_g is the graph’s adjacency matrix, which encodes the connectivity between nodes. $W^{(k)}$ is the weight matrix of the k -th layer, which is used to learn a linear transformation of the node features and σ is a non-linear activation function, such as ReLU or sigmoid.

After the message passing, the aggregated information vector ζ'_i for each agent i , for a GraphNet having k hidden layers, is given as,

$$\zeta'_i = H^{(k)} = \sigma(A_g H^{(k-1)} W^{(k-1)}). \quad (7)$$

The aggregated information vector ζ' is now passed on to the actor-critic network MLP. The actor-network makes decisions for the agent, and a separate critic network evaluates the actor’s actions to provide feedback, allowing the actor to improve its decision-making over time.

C. Multi-agent actor-critic method using PPO

The decentralized actors in the multi-agent PPO take in the aggregated information vector ζ'_i and generate the corresponding action probability distribution π . The action space consists of five discrete options: {up, down, left, right, stay}, representing decisions to move in one of the four cardinal directions or to remain in the current location.

The centralized critic estimates the environment’s value function to influence the individual actor’s policy (Figure 3). The shared reward for all agents is defined in Equation (2).

For a defined episode length T , the agent interacts with the environment to generate and collect the trajectory values in the form of states, actions, and rewards $\{s_i, a_i, r_i\}$. The stored values are then sampled iteratively to update the action probabilities and to fit the value function through back-propagation.

Let θ_1 be the actor trainable parameter and θ_2 be the critic trainable parameter. Discounted return measures the long-term value of a sequence of actions. The discounted return is given as $G(t; \theta_1) = \sum_{k=0}^T \gamma^k r(t+k+1; \theta_1)$, where $\gamma \in [0, 1)$ is the discount factor and T is the episode time horizon. The Value function $V(s_t^i; \theta_2)$ represents the expected long-term reward that an agent i can expect to receive if it starts in that state s at time t . It is updated as the agent interacts with the environment and learns from its experiences. The value function estimate, which is defined as $V(s_t^i; \theta_2) = \mathbb{E}[G(t)|s_t^i]$, is provided by the critic network. The advantage estimate function \hat{A}_i is a measure of how much better a particular action is compared to the average action taken by the current policy. It is defined as the difference between the discounted return and the state value estimate given by

$$\hat{A}_t^i(\theta_1, \theta_2) = G(t; \theta_1) - V(s_t^i; \theta_2). \quad (8)$$

PPO uses the advantage function to adjust the probability of selecting an action to make the policy more likely to take actions with a higher advantage. This helps ensure that the policy makes the most efficient use of its resources and maximizes the expected reward over time [14]. The modified multi-agent PPO objective function to be minimized in the GALOPP network is given as,

$$L(\theta_1, \theta_2) = \frac{1}{m} \sum_m \left(\frac{1}{N} \sum_{i=1}^N (L_i^{CLIP}(\theta_1, \theta_2)) \right), \quad (9)$$

where N is the total number of agents and m is the mini-batch size, and $L_i^{CLIP}(\theta_1, \theta_2)$ refers to the clipped surrogate

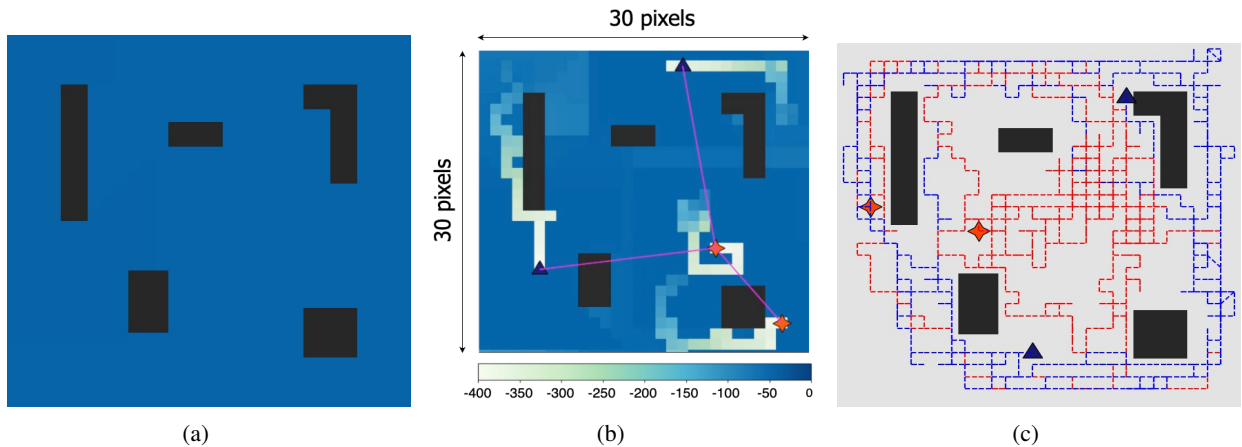


Fig. 5: (a) Outline of the open map. (b) The agents cannot move into black pixels, while the non-black regions need to be persistently monitored. As the anchor agents (red stars) and auxiliary agents (dark blue triangles) monitor, their trajectory is shown as the fading white trails for the last 30 steps. The communication range between the agents is shown in red lines. (c) The trajectories of the anchor and auxiliary agents while monitoring are shown by the red and blue lines, respectively.

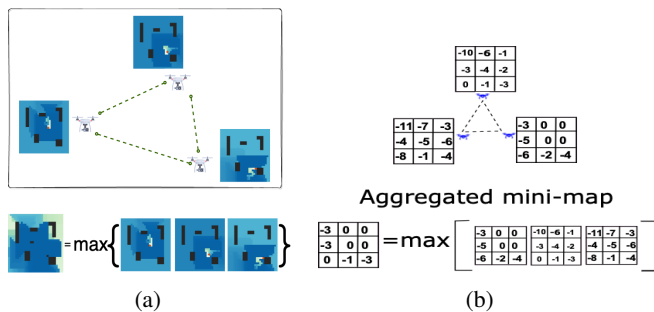


Fig. 6: (a) Illustration of decentralized map-sharing among agents in persistent monitoring. (b) Overview of how agents within communicable range of one another update their global maps in a decentralized setting. The resultant global map is generated by taking the element-wise maximum value from the individual global maps of the agents.

objective function [14] defined as

$$L_i^{CLIP}(\theta_1, \theta_2) = \hat{\mathbb{E}}_t[\min(r_t(\theta_1)\hat{A}_t^i(\theta_1, \theta_2), \text{clip}(r_t(\theta_1), 1 - \epsilon, 1 + \epsilon)\hat{A}_t^i(\theta_1, \theta_2))], \quad (10)$$

where $r_t(\theta_1) = \pi_{\theta_1} / \pi_{\theta_1^{\text{old}}}$ is the current policy's (π_{θ_1}) action probability ratio to the previous policy distribution $\pi_{\theta_1^{\text{old}}}$. The *clip* function clips the probability ratio $r_t(\theta_1)$ to the trust-region interval $[1 - \epsilon, 1 + \epsilon]$.

GALOPP is trained end to end by minimizing the modified PPO objective function using the trajectory values collected from the interactions with the environment. GALOPP minimizes the multi-agent PPO objective function to train the network. The algorithm updates the action probabilities and fits the value function through back-propagation. This allows the model to learn from experience and improve its performance over time.

V. EXPERIMENTS AND ANALYSIS

We evaluate the performance of GALOPP on an open map environment as shown in Figure 5. The open map has an area of 30×30 sq. units, where 5 obstacles having random geometry are placed. The agents have a sensing range of $\ell = 7$ in the 2D environment. We use the accumulated reward metric to evaluate the performance. The total reward at time t is defined as $\mathcal{R}(t) = \sum_{\alpha, \beta} R_{\alpha\beta}(t)$. The grid cells' penalties are updated with a decay rate of $\Delta_{\alpha\beta} = 1, \forall(\alpha, \beta)$. A cell's maximum negative reward is $R_{max} = 400$. The simulation parameters used in the experiments are detailed in Table I.

A. Model

GALOPP was trained and tested using Python 3.6 on a workstation with Ubuntu 20.04 L.T.S operating system, with an Intel(R) Core(TM) i9 CPU and an NVIDIA GeForce RTX 3090 GPU (running on CUDA 11.1 drivers). The neural networks were written and trained using PyTorch 1.8 and dgl-cu111 (deep graph library). We now provide details of the various parameters used in the model. The GALOPP architecture consists of 4 deep neural networks: ConvNet, GraphNet, Actor MLP, and Critic MLP, as shown in Figure 3. The details of this architecture are given in Table II.

1) *Embedding generator (ConvNet)*: This convolutional neural network takes a 2-channelled 15×15 image (local map and mini-map) as the input and generates a 32-dimensional feature vector. We then append a 6-dimensional state vector to this feature vector (positional mean and covariance) to form a 38-dimensional feature vector that acts as the embedding for the graph convolution network. The state vector is derived by flattening the agent's covariance matrix Σ_t and appending it to the position vector μ_t .

2) *Graph convolution network (GraphNet)*: The embeddings generated by the embedding generator are passed through a single-layered feed-forward graph convolution network to generate the embeddings for the actor networks of the individual agents.

Parameter	Value
Decay Rate ($\Delta_{\alpha\beta}$)	1
Maximum negative reward (R_{max})	400
Length of episode (T)	1000
Agent visibility range (ℓ)	7
Local map and Mini map size ($g \times g$)	15×15
O_t (covariance matrix for error in IMU suite)	$\begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}$
Q_t (covariance matrix for uncertainty in sensors)	$\begin{bmatrix} 1e-4 & 0 \\ 0 & 1e-4 \end{bmatrix}$

TABLE I: Simulation Parameters for GALOPP

ConvNet
ConvLayer1 (in-channels=2, out-channels=16, kernel-size=8, stride=4, padding=(1, 1)), ReLU activation
ConvLayer2 (in-channels=16, out-channels=32, kernel-size=4, stride=2, padding=(1, 1)), ReLU activation
ConvLayer3 (in-channels=32, out-channels=32, kernel-size=3, stride=1, padding=(1, 1)), ReLU activation, Flatten
GraphNet
GCNLayer(in-features=38, out-features=38)
Actor MLP
LinearLayer1 (in-features=38, out-features=500), ReLU activation
LinearLayer2 (in-features=500, out-features=256), ReLU activation
LinearLayer3 (in-features=256, out-features=5), SoftMax
Critic MLP
LinearLayer1 (in-features=38, out-features=500), ReLU activation
LinearLayer2 (in-features=500, out-features=256), ReLU activation
LinearLayer3 (in-features=256, out-features=1)

TABLE II: Parameters for the neural networks

3) *Actor MLP*: The actor takes the embeddings generated by the ConvNet and the aggregated information vector from the GraphNet network as the input and generates the probability distribution for the available actions.

4) *Critic MLP*: The critic network takes the embeddings generated by the ConvNet for each agent and returns the state-value estimate for the current state.

5) *Training*: The training is carried out for 30000 episodes where each episode is of length $T = 1000$ time steps. The agents are initialized randomly in the environment for every training episode but are localized during initialization.

The GALOPP architecture input at time t is the image representing the state of the grid G , which is resized to an image of the dimension 15×15 using OpenCV's INTER_AREA interpolation method and concatenated with the local visibility map of the agent forming a 2-channeled image of dimension 15×15 . The action space has five actions: up, down, left, right, and stay. Each action enables the agent to move by one pixel, respectively.

6) *Evaluation*: For testing the learned policies, we evaluate it for 100 episodes, each episode for $T = 1000$ time steps, in their respective environments. The reward for test

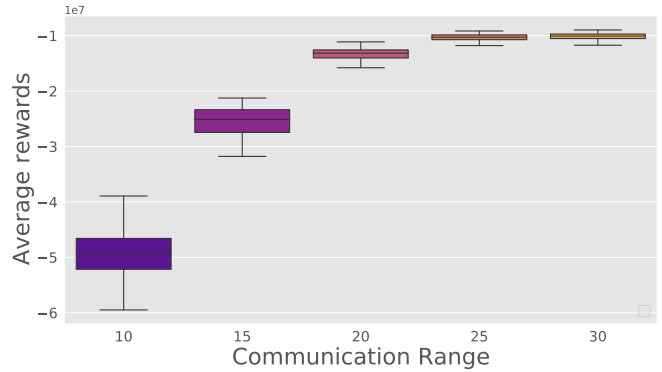


Fig. 7: Comparison of the average reward on increasing the communication range of the agents in the open-map environment.

episode τ is denoted by $\mathcal{R}_{ep}^\tau = \sum_{t=1}^T \mathcal{R}(t)$ and the final reward \mathcal{R}_{avg} after $n = 100$ episodes are calculated as $\mathcal{R}_{avg} = \frac{1}{n} \sum_{\tau=1}^n \mathcal{R}_{ep}^\tau$. The \mathcal{R}_{avg} is used to evaluate the model's performance. Next, we will evaluate the performance of the GALOPP under different parameters.

B. Effect of increase in communication range

With an increase in communication, the agents will be able to communicate as well as localize better while reaching various locations in the environment. A lower communication range can make agents close to each other, and hence, the agents are unable to explore and cover different regions, making an ineffective strategy. We consider a system comprising 2 anchor agents and 2 auxiliary agents, and vary the communication range from $\rho = 10$ units to $\rho = 30$ units, with an increment of 5 units. We evaluate the performance of GALOPP under different communication ranges as shown in Fig. 7.

From the figure, we can see that with a reduced communication range of 10 and 15, the agents are unable to monitor the region properly, hence resulting in higher negative rewards. As we increase the communication range to 20, the performance improves as the agents can communicate better while maintaining localization accuracy. However, by increasing the communication range higher than 20, there is a marginal improvement in the performance at the cost of a higher communication range. These results are intuitive. However, they provide insight into the selection of the communication range for the rest of the simulations. Based on these results, we consider $\rho = 20$ for the rest of the analysis.

C. Effect of varying sensing range

The size of the local map is dependent on the sensing range ℓ , which we measure in terms of the number of cells that can be observed. As the sensing range ℓ increases, the number of observed cells $g \times g$ also increases, where $g = 2\ell + 1$, resulting in a decrease in penalties. Intuitively, with an increase in sensing range, the reward improves, which can be seen in Figure 8. The difference in performance

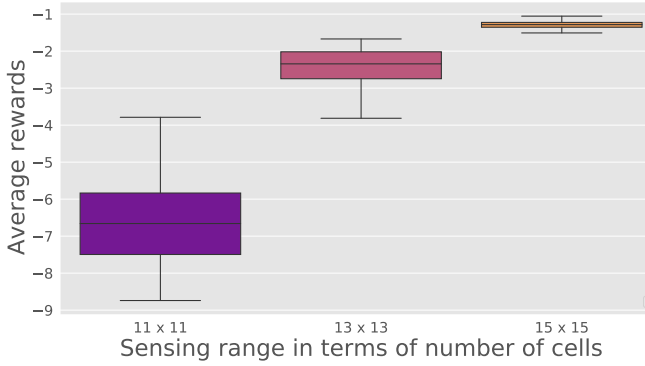


Fig. 8: Comparison of the average reward of the model on decreasing the local sensing map range. The local map is the agents’ visibility range in the environment.

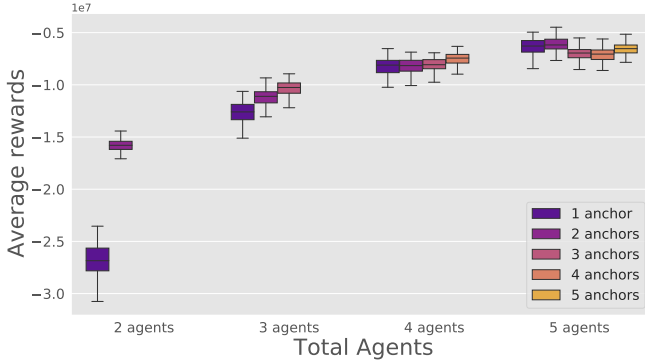


Fig. 9: Effect of increasing the total number of agents in the environment. For a given number of agents, we effect of increasing the number of anchor agents $k \leq N$ for N agents in the environment.

between $\ell = 5$ and $\ell = 6$ is significant; however, the performance improvement is lower when we further increase the sensing range to $\ell = 7$. Based on these trends, if we further increase the sensing range, the improvement will be marginal. Hence, we consider a sensing range of $\ell = 7$ for the rest of the simulation. Note that during this evaluation, we use a communication range of $\rho = 20$, as fixed from the previous analysis.

D. Effect of an increasing number of agents and varying anchor-auxiliary ratio

The ability to monitor adequately in the environment depends on the number of agents present in the environment and also the ratio of the anchor to auxiliary agents. To understand this effect, we carry out simulations, varying the number of agents from 2 to 5. For a given agent, we vary the number of anchors to understand the performance to cost benefits associated with a higher number of anchors. Figure 9 shows the model performance for a varying number of agents in the environment. First, let’s consider the effect of an increase in the number of agents with a single anchor. From the figure, we can see that with an increase in the number of agents, the coverage is higher and hence improvement in the average rewards. However, as we increase the number from

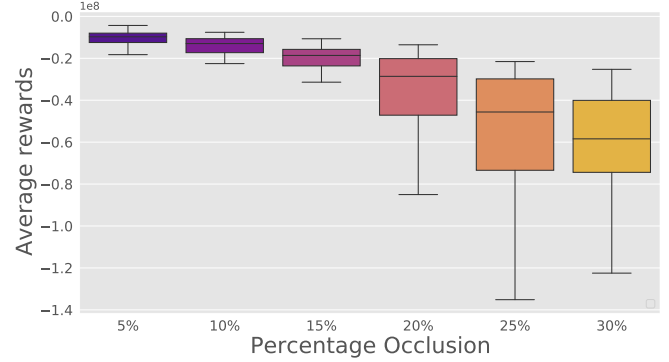


Fig. 10: Comparison of the average reward on increasing the percentage obstruction in the environment by increasing the number of obstacle blocks.

4 to 5, the improvement is marginal because four agents are sufficient to cover the region, and hence, increasing more agents does not increase the rewards significantly.

For a given number of agents, let us now analyze the effect of a number of anchor agents. For 2 agents, with both being anchors enables the agents to cover better, and since these two agents have high accuracy, they can work independently, thus improving the performance of one anchor. When we increase the number of anchors for 3, 4, and 5 agent cases, we can see that increasing the number of anchors shows only a marginal improvement. Hence, we can obtain good coverage accuracy with a lower number of anchor agents while ensuring there are 2 or more auxiliary agents. With a lower number of anchors, the deployment cost can be reduced significantly.

E. Effect of increasing obstruction in the environment

The model should have the robustness to be able to perform well under different percentage obstructions in the environment. However, as the percentage of obstructions increases, the difficulty in monitoring also increases. In order to validate this hypothesis, we perform simulations on varying obstacle percentages in the environment. For each episode, the obstacles for a given percentage are randomly generated and placed.

Figure 10 shows the performance of GALOPP for varying percentage obstruction. From the figure, we can see that when the obstruction is less (5-15%), the GALOPP model is able to learn to change the paths so that the rewards are maximized. However, with further increases in obstacle density (20-30%), learning becomes difficult due to environmental constraints and, hence reduction in performance. When we look at the percentage of disconnections that happen due to environmental changes, for 5-15% obstacle density, the disconnections are less than 10%. However, with an increase in the obstacle density, the motion constraints for the agents also increase. Due to this, the agents are unable to explore remote regions, resulting in reduced performance as shown in Figure 10. Because the agents are unable to disconnect and explore, they remain connected, resulting in a lower percentage of disconnection time.

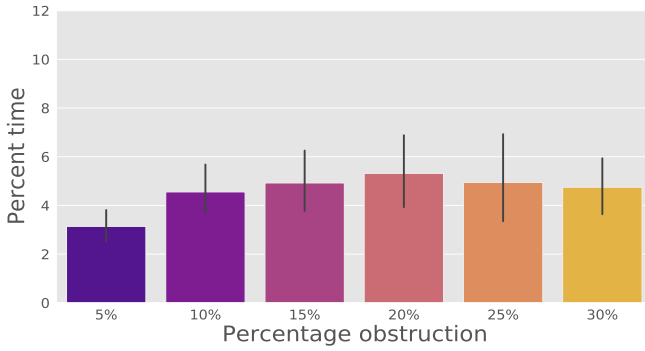


Fig. 11: Comparison of the percent time of disconnection for auxiliary agents on increasing the percentage occlusion in the environment.

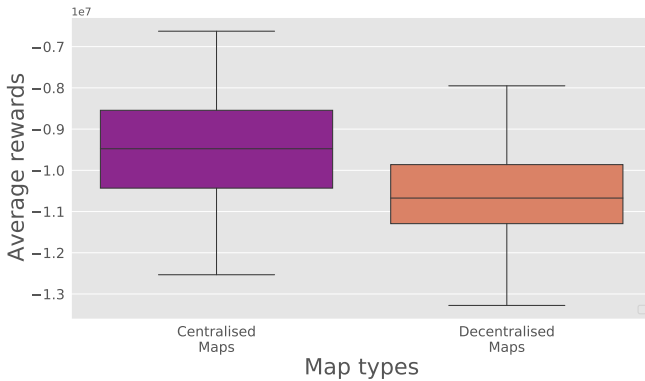


Fig. 12: Comparison between centralized and decentralized execution

F. Comparison between centralized maps vs. decentralized maps

In GALOPP, agents are trained using a decentralized mini-map, where each agent maintains a separate copy of the global map that was updated when agents were within a communicable range. We compare the performance of the decentralized global map approach to a centralized approach, where a shared global map is maintained among all agents. To accomplish this comparison, agents within the communication range of each other compared and aggregated the global map at each time step by taking the element-wise maximum for each grid cell in the environment, as shown in Figure 6. In order to know the difference in performance between centralized map sharing and decentralized map sharing, simulations were carried out, and Figure 12 shows the performance difference. The simulations setting for the comparison are two anchor agents and two auxiliary agents with a sensing range of 7 cells and a communication range of 20.

From the figure, we can see that the centralized map model is performing marginally better than the decentralized map model, but statistically, both strategies are performing similarly. The result shows that using decentralized maps is a good alternative to centralized maps. This suggests that the decentralized approach in GALOPP can achieve

similar performance to a centralized approach while still providing the benefits of decentralization in maintaining its local observation.

G. Comparison between GALOPP and non-RL baselines

Due to the localization constraints in the persistent monitoring problem, the motion of both anchor and auxiliary agents becomes coupled. This complexity makes it highly challenging to devise deterministic motion strategies for these heterogeneous agents. Given the unique coupled objective of coverage and localization in our approach, existing literature has limitations in simultaneously addressing both aspects, which in turn complicates the identification of suitable baselines. Since there is limited prior research on persistent monitoring with constraints on communication, localization, and sensing range, we compare our approach with heuristic-based algorithms custom-designed for persistent monitoring.

As a result, we assess the performance of our model against four custom-designed non-reinforcement learning baselines: random search (RS), random search with ensured communication (RSEC), greedy search (GS), and lawn-mower area sweep (LMAS).

1) *Random Search (RS)*: In the RS method, agents make decisions independently at each time step by randomly selecting an action (stay, up, down, right, left). This approach does not require any prior knowledge of the problem domain or any model of the system dynamics. Because of random decisions, communication may break, resulting in lower performance.

2) *Random Search with Ensured Communication (RSEC)*: RSEC is an extension of the RS method, in which each agent randomly selects an action while ensuring that no auxiliary agent becomes unlocalized. In other words, the RSEC approach guarantees that all agents remain localized at all times. If an action is selected that would cause an agent or another auxiliary agent to become unlocalized, the agent randomly selects another action from the remaining action space until a suitable action is found.

3) *Greedy Search (GS)*: In GS, agents act independently and greedily. Assume that agent i is in cell (α, β) , and we define \mathcal{N}_i as the set of neighboring cells that agent i can reach in one-time step (that is, all the cells when $\ell = 1$). The agent i selects a cell that has maximum negative reward, without considering localization constraints. If all the grid cells in \mathcal{N}_i have the same reward, then agent i chooses a random action.

4) *Area partition with Lawn Mower Sweep (LMAS)*: The LMAS strategy [18], [38] begins with partitioning the area into cells (or sections) equal to the number of agents deployed and placing each agent in a chosen starting position within one of these sections. The agent is then programmed to follow a specific lawn-mower movement pattern within each section, typically involving a back-and-forth motion. The agent repeats this pattern within each section until it has covered the entire area, ensuring complete coverage. Each agent makes use of its onboard local map to avoid obstacles by changing its trajectory in the presence of obstacles.

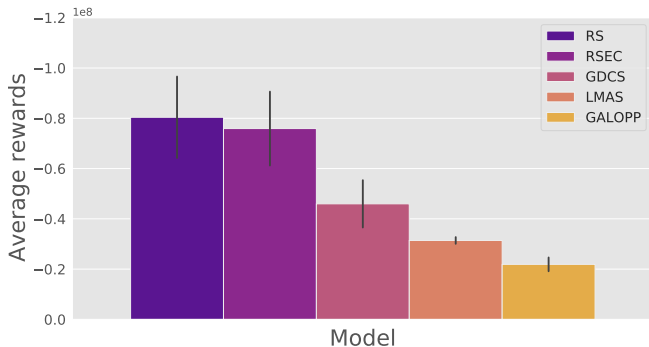


Fig. 13: Performance Comparison of GALOPP with Heuristic Baselines: Random Search (RS), Random Search with Ensured Communication (RSEC), Greedy Search (GS), and Lawn Mower Area Sweep (LMAS)

We carried out 100 simulations for each non-RL baseline strategy and Figure 13 shows the performance comparison between the baseline strategies and GALOPP. From the figure, we can see that the GALOPP consistently outperforms the above-defined baseline strategies. This is attributed to GALOPP’s ability to explicitly account for localization and connectivity constraints in its decision-making process. Among the baselines, Random Search (RS) exhibits the poorest performance, as it relies on random actions without considering any context. Random Search with Ensured Communication (RSEC) improves upon RS by enforcing communication and localization, resulting in enhanced performance. Greedy Search (GS) leads to sub-optimal policies as each agent acts greedily at each timestep independently. Lawn Mower Area Sweep (LMAS) performs better compared to GS, but it does not explicitly incorporate localization considerations and is influenced by the specific geometry of the environment being monitored. GALOPP consistently outperforms all the mentioned baselines.

H. Evaluation in other environments

In order to test the ability of GALOPP to perform in other types of complex environments, we evaluate its performance in two-room and four-room environments, as shown in Figure 14a and 14c, respectively.

For the two-room map, the agents learn to maintain contact with each other by spreading across two rooms and the corridor. In the 2-room map, we notice that our algorithm ends up with the agents in a formation where two of them position themselves in the two rooms while one monitors the corridor. This can be seen in Figure 14a, where the faded cells show the trajectory followed by the agents for the last 30 steps. Figure 14b shows the areas where each agent was present. From this, we can see that the anchor was in the middle region while the two auxiliary agents monitored the two rooms. The anchor agent moves around to maximize the rewards, while the auxiliary agents move in the two rooms. In fact, this is the best combination for the agents, and they learn quickly.

In the four-room map, GALOPP learns a policy in which

each of the four agents is responsible for monitoring separate rooms while intermittently monitoring the central corridor region, as shown in Figures 14c and 14d. The anchor agents are positioned to monitor two cells and the central area, while the auxiliary agents are responsible for monitoring the two rooms.

Our results show that GALOPP is capable of adapting to complex environments and learning effective policies for multi-agent coordination. The ability of the agents to maintain contact with each other and cover all areas of the environment is crucial for the successful completion of tasks, and GALOPP demonstrates its ability to achieve this.

VI. HARDWARE IMPLEMENTATION

We implement GALOPP on a real-time hardware setup for proof-of-concept purposes. We use multiple BitCraze Crazyflie 2.1 [39] nanocopters as agents. The experimental setup consists of four SteamVR Base stations [40] and Lighthouse Positioning System [41] to track the location of the vehicles within a $3.5m \times 3.0m \times 2.0m$ arena. The agents communicate with a companion computer (running on Ubuntu 20.04 with an AMD Ryzen 9 5950x with a base clock speed of 3.4 GHz) via a Crazyradio telemetry module, where the trained GALOPP model was executed. In the experiment, we consider the environment as shown in Fig. 15a with 2 auxiliary agents and 1 anchor agent. The companion computer receives the position of each CrazyFlie as input via the corresponding rostopics from the CrazySwarm ROS package [42] [43]. The respective agents then execute the actions computed by the actor networks. To avoid inaccuracies in tracking the CrazyFlies caused by physical obstacles obstructing the infrared laser beams from the Base stations, we opt to simulate the obstacle boundaries. The model policy implemented in the simulation ensures that the agents never collide with any obstacle.

The video of the hardware implementation can be seen in [44]. Figure 15(a) shows the snapshot of the simulated environment along with the agent positions (anchor and auxiliary), current coverage, and the position of the obstacle. We then implement the same scenario with virtual obstacles through the hardware, where the model sends the control signals to the vehicles, as shown in Figure 15(b). In Figure 15(c), we can see that the agent trajectories are covering all the regions and hence achieving persistent monitoring.

VII. CONCLUSION AND FUTURE WORK

This work developed a MARL algorithm with a graph-based connectivity approach – GALOPP for persistent monitoring a bounded region, taking the communication, sensing, and localization constraints into account via graph connectivity. The experiments show that the agents using GALOPP can outperform four custom baseline strategies for persistent area coverage while accounting for the connectivity bounds. We also establish the robustness of our approach by varying the sensing map, the effect of obstacle occlusion by increasing the percent amount of obstacle, and by scaling the number of anchor agents in the system.

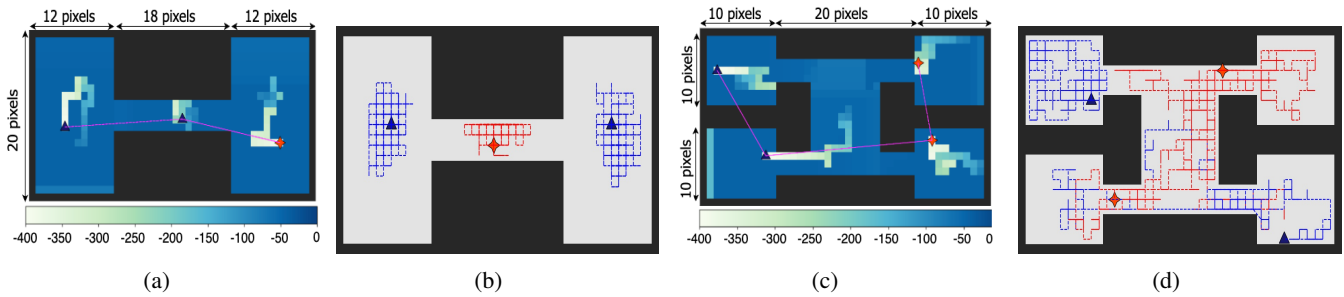


Fig. 14: Visualization of maps: (a) Illustrates a 2-room map, and (c) Illustrates a 4-room map. The agents cannot move into black pixels, while the non-black regions need to be persistently monitored. As the anchor agents (red stars) and auxiliary agents (dark blue triangles) monitor, their trajectory is shown as the fading white trails for the last 30 steps. The communication range between the agents is shown in red lines. Fig (b) and Fig (d) display the trajectories of the anchor and auxiliary agents while monitoring for the 2-room and 4-room maps, respectively.

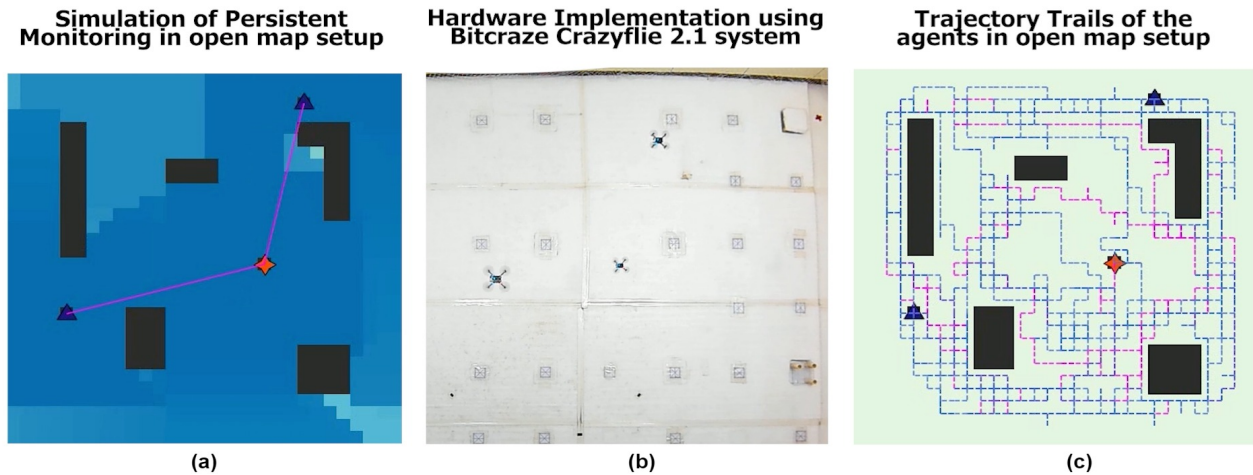


Fig. 15: Snapshot of the video for the hardware implementation of vehicles using one anchor and two auxiliary agents. (a) A rendered simulation snapshot of the monitoring task. (b) Real-time decision-making being performed by the trained GALOPP network model. (c) The trajectory trails of the previous timesteps that the agent took in the monitoring task.

It was seen that increasing the number of anchor agents improves the performance, but beyond a certain value, there are diminishing returns on the rewards obtained. Based on power and resource constraints, one can select a subset of anchor agents to achieve persistent surveillance effectively.

Although our experiments demonstrate that GALOPP surpasses the baseline strategies, future work could investigate the algorithm’s scalability as the number of agents significantly increases. Exploring optimal values for the decay rate and maximum negative reward, with a focus on increasing monitoring efficiency, presents a promising area for further research. Additionally, the algorithm’s suitability for diverse sensor types, such as cameras or LIDAR sensors, could be explored to improve agents’ situational awareness. Further research on the impact of different types of obstacles, including moving obstacles, on the algorithm’s performance would also be insightful. While the proposed algorithm targets heterogeneous agents in the persistent monitoring problem, future research can investigate its generalizability to other monitoring problems, such as target tracking or

environmental monitoring. This work provides a foundation for future investigations of GALOPP’s performance and its potential applications in various monitoring scenarios.

REFERENCES

- [1] J. Yu, S. Karaman, and D. Rus, “Persistent monitoring of events with stochastic arrivals at multiple stations,” *IEEE Transactions on Robotics*, vol. 31, no. 3, pp. 521–535, 2015.
- [2] S. L. Smith, M. Schwager, and D. Rus, “Persistent monitoring of changing environments using a robot with limited range sensing,” in *IEEE International Conference on Robotics and Automation*, 2011, pp. 5448–5455.
- [3] S. K. K. Hari, S. Rathinam, S. Darbha, K. Kalyanam, S. G. Manyam, and D. Casbeer, “The generalized persistent monitoring problem,” in *American Control Conference*, 2019, pp. 2783–2788.
- [4] X. Lin and C. G. Cassandras, “An optimal control approach to the multi-agent persistent monitoring problem in two-dimensional spaces,” *IEEE Transactions on Automatic Control*, vol. 60, no. 6, pp. 1659–1664, 2014.
- [5] S. K. K. Hari, S. Rathinam, S. Darbha, K. Kalyanam, S. G. Manyam, and D. Casbeer, “Optimal uav route planning for persistent monitoring missions,” *IEEE Transactions on Robotics*, vol. 37, no. 2, 2021.
- [6] T. Wang, P. Huang, and G. Dong, “Cooperative persistent surveillance on a road network by multi-ugvs with detection ability,” *IEEE Transactions on Industrial Electronics*, vol. 69, no. 11, pp. 11 468–11 478, 2021.

- [7] V. Mersheeva and G. Friedrich, "Multi-uav monitoring with priorities and limited energy resources," in *Proceedings of the International Conference on Automated Planning and Scheduling*, vol. 25, 2015, pp. 347–355.
- [8] Y.-W. Wang, Y.-W. Wei, X.-K. Liu, N. Zhou, and C. G. Cassandras, "Optimal persistent monitoring using second-order agents with physical constraints," *IEEE Transactions on Automatic Control*, vol. 64, no. 8, pp. 3239–3252, 2018.
- [9] E. Arribas, V. Cholvi, and V. Mancuso, "Optimizing uav resupply scheduling for heterogeneous and persistent aerial service," *IEEE Transactions on Robotics*, 2023.
- [10] J. Zhu and S. S. Kia, "Cooperative localization under limited connectivity," *IEEE Transactions on Robotics*, vol. 35, no. 6, pp. 1523–1530, 2019.
- [11] J. Liu, J. Pu, L. Sun, and Y. Zhang, "Multi-robot cooperative localization with range-only measurement by uwb," in *Chinese Automation Congress*. IEEE, 2018, pp. 2809–2813.
- [12] R. Sharma, R. W. Beard, C. N. Taylor, and S. Quebe, "Graph-based observability analysis of bearing-only cooperative localization," *IEEE Transactions on Robotics*, vol. 28, no. 2, pp. 522–529, 2011.
- [13] F. Klaesson, P. Nilsson, T. S. Vaquero, S. Tepsuporn, A. D. Ames, and R. M. Murray, "Planning and optimization for multi-robot planetary cave exploration under intermittent connectivity constraints," in *ICAPS Workshop on Planning and Robotics*, 2020.
- [14] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.
- [15] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *International Conference on Learning Representations*, 2017.
- [16] J. O’rouke *et al.*, *Art gallery theorems and algorithms*. Oxford University Press Oxford, 1987, vol. 57.
- [17] P. Tokekar and V. Kumar, "Visibility-based persistent monitoring with robot teams," in *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2015, pp. 3387–3394.
- [18] H. Choset, "Coverage for robotics—a survey of recent results," *Annals of mathematics and artificial intelligence*, vol. 31, pp. 113–126, 2001.
- [19] E. Galceran and M. Carreras, "A survey on coverage path planning for robotics," *Robotics and Autonomous systems*, vol. 61, no. 12, pp. 1258–1276, 2013.
- [20] X. Tan, "Fast computation of shortest watchman routes in simple polygons," *Information Processing Letters*, vol. 77, no. 1, pp. 27–33, 2001.
- [21] P. H. Washington and M. Schwager, "Reduced state value iteration for multi-drone persistent surveillance with charging constraints," in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2021, pp. 6390–6397.
- [22] P. Maini, P. Tokekar, and P. B. Sujit, "Visibility-based persistent monitoring of piecewise linear features on a terrain using multiple aerial and ground robots," *IEEE Transactions on Automation Science and Engineering*, vol. 18, no. 4, pp. 1692–1704, 2021.
- [23] S. Omidshafiei, J. Papis, C. Amato, J. P. How, and J. Vian, "Deep decentralized multi-task multi-agent reinforcement learning under partial observability," in *International Conference on Machine Learning*, Sydney, Australia, 2017, pp. 2681–2690.
- [24] D. Maravall, J. de Lope, and R. Domínguez, "Coordination of communication in robot teams by reinforcement learning," *Robotics and Autonomous Systems*, vol. 61, no. 7, pp. 661–666, 2013.
- [25] Q. Li, F. Gama, A. Ribeiro, and A. Prorok, "Graph neural networks for decentralized multi-robot path planning," *IEEE International Conference on Intelligent Robots and Systems*, 2019.
- [26] R. Shah, Y. Jiang, J. Hart, and P. Stone, "Deep r-learning for continual area sweeping," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Las Vegas, 2020, pp. 5542–5547.
- [27] Q. Li, W. Lin, Z. Liu, and A. Prorok, "Message-aware graph attention networks for large-scale multi-robot path planning," *IEEE Robotics and Automation Letters*, vol. 6, no. 3, pp. 5533–5540, 2021.
- [28] B. Wang, Z. Liu, Q. Li, and A. Prorok, "Mobile robot path planning in dynamic environments through globally guided reinforcement learning," *IEEE Robotics and Automation Letters*, vol. 5, no. 4, pp. 6932–6939, 2020.
- [29] J. Blumenkamp and A. Prorok, "The emergence of adversarial communication in multi-agent reinforcement learning," *Conference on Robot Learning*, Cambridge, MA, USA, 2020.
- [30] J. Chen, A. Baskaran, Z. Zhang, and P. Tokekar, "Multi-agent reinforcement learning for visibility-based persistent monitoring," in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2021, pp. 2563–2570.
- [31] F. Klaesson, P. Nilsson, A. D. Ames, and R. M. Murray, "Intermittent connectivity for exploration in communication-constrained multi-agent systems," in *ACM/IEEE International Conference on Cyber-Physical Systems*, 2020, pp. 196–205.
- [32] R. Khodayi-mehr, Y. Kantaros, and M. M. Zavlanos, "Distributed state estimation using intermittently connected robot networks," *IEEE Transactions on Robotics*, vol. 35, no. 3, pp. 709–724, 2019.
- [33] S. Thrun, "Probabilistic robotics," *Communications of the ACM*, vol. 45, no. 3, pp. 52–57, 2002.
- [34] R. Lowe, Y. Wu, A. Tamar, J. Harb, P. Abbeel, and I. Mordatch, "Multi-agent actor-critic for mixed cooperative-competitive environments," in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, Long Beach, California, USA, 2017, p. 6382–6393.
- [35] C. Yu, A. Velu, E. Vinitzky, J. Gao, Y. Wang, A. Bayen, and Y. Wu, "The surprising effectiveness of ppo in cooperative multi-agent games," *Advances in Neural Information Processing Systems*, vol. 35, pp. 24 611–24 624, 2022.
- [36] Y. LeCun, Y. Bengio *et al.*, "Convolutional networks for images, speech, and time series," *The handbook of brain theory and neural networks*, vol. 3361, no. 10, p. 1995, 1995.
- [37] J. Chen, A. Baskaran, Z. Zhang, and P. Tokekar, "Multi-agent reinforcement learning for persistent monitoring," *arXiv preprint arXiv:2011.01129*, 2020.
- [38] H. Choset, K. M. Lynch, S. Hutchinson, G. A. Kantor, and W. Burgard, *Principles of robot motion: theory, algorithms, and implementations*. MIT press, 2005.
- [39] A. BITCRAZE, "Crazyflie 2.1." [Online]. Available: <https://www.bitcraze.io/products/crazyflie-2-1/>
- [40] H. VIVE, "Basestation." [Online]. Available: <https://www.vive.com/sea/accessory/base-station2/>
- [41] A. BITCRAZE, "Lighthouse positioning system." [Online]. Available: <https://www.bitcraze.io/documentation/tutorials/getting-started-with-lighthouse/>
- [42] J. A. Preiss, W. Honig, G. S. Sukhatme, and N. Ayanian, "Crazyswarm: A large nano-quadcopter swarm," in *2017 IEEE International Conference on Robotics and Automation (ICRA)*, 2017, pp. 3299–3304.
- [43] M. Quigley, K. Conley, B. Gerkey, J. Faust, T. Foote, J. Leibs, R. Wheeler, A. Y. Ng *et al.*, "Ros: an open-source robot operating system," in *ICRA workshop on open source software*, vol. 3, no. 3.2. Kobe, Japan, 2009, p. 5.
- [44] "Proof-of-concept hardware experiment," https://moonlab.iiserb.ac.in/research_page/galopp.html.