

Multi-Robot Coverage and Exploration using Spatial Graph Neural Networks

Ekaterina Tolstaya¹, James Paulos², Vijay Kumar^{1,2}, Alejandro Ribeiro¹

Abstract—The multi-robot coverage problem is an essential building block for systems that perform tasks like inspection, exploration, or search and rescue. We discretize the coverage problem to induce a spatial graph of locations and represent robots as nodes in the graph. Then, we train a Graph Neural Network controller that leverages the spatial equivariance of the task to imitate an expert open-loop routing solution. This approach generalizes well to much larger maps and larger teams that are intractable for the expert. In particular, the model generalizes effectively to a simulation of ten quadrotors and dozens of buildings in an urban setting. We also demonstrate the GNN controller can surpass planning-based approaches in an exploration task.

I. INTRODUCTION

Large scale swarms of robots could be deployed to provide on-demand wireless networks [1], perform rapid environmental mapping [2], [3], track targets [4], search after natural disasters [5], [6], or enable sensor coverage in communication-denied environments [7]. At moderate scales, it may be possible to centralize the entire team’s information and control in one agent, but practical deployments of very large teams require distributed execution and scalable algorithms. In particular, global planning-based approaches suffer from an exponential increase in complexity as the number of robots and the environment size increases. This motivates the use of heuristics in general and, as we advocate in this paper, the use of learned heuristics.

Graph neural networks (GNNs) have been used to generate heuristic solutions to a variety of multi-robot problems, such as path planning [8], [9], [10], exploration [11], and perimeter defense [12]. In large teams, we can take advantage of graph equivariance to design abstractions that further speed up learning. In this work, we focus on the problem of coverage, in which a robot team must visit a set of locations in an environment [13]. To leverage recent advances in graph neural networks, we encode the task as a graph: the known map locations and team members are graph nodes, and allowed moves are graph edges. This approach allows us to abstract away the global agent and obstacle locations and to represent all elements of the problem in a single spatial graph with only local connectivity. Most importantly, the

spatial graph representation can describe tasks with complex spatial constraints, unlike past work that applies GNNs to homogeneous inter-robot communication graphs [14], [15].

A moderate-size coverage task with dozens of goals and fewer than ten agents can be solved with existing approaches when posed as a vehicle routing problem [16], [17]. We collect a dataset of trajectories generated using the centralized expert solution and use behavior cloning to train a graph neural network controller to imitate the expert solution. This learned heuristic can then generalize to previously unseen coverage scenarios with more agents and larger maps. We also show generalization to a scenario with simulated quadrotors that must traverse an environment with thousands of waypoints. Furthermore, we apply this approach to an exploration task, which is similar to the coverage task but the graph of waypoints is revealed to the robot team during mission execution. We demonstrate that by imitating an omniscient expert, our learned controller can outperform both a greedy controller and a receding horizon planner.

For the applications of coverage and exploration, information about distant points of interest may be necessary for computing the next position for each robot. To enable the learned controllers to use distant information, we build graph neural networks with a larger number of graph operation layers, up to 19 layers. The number of graph operation layers determines the distance along which information can travel from one node to another along the edges in the graph, the *receptive field* of the architecture. Other applications of GNNs to robot teams typically use receptive fields of 2 to 4 in conjunction with dense adjacency matrices [11], [15], [14], [18], with the exception of the shortest path demonstration in [9] with a receptive field of 10. A mean aggregation operation helps stabilize training of GNNs with larger receptive fields. Furthermore, we use a sparse representation of the local connectivity of the graph that allows our approach to scale to larger maps and teams than [18]. Our contributions can be summarized as follows:

- 1) An approach to encoding continuous space multi-robot coverage and exploration problem data as a discrete spatial graph in which robots and points of interest are nodes and allowed moves are edges.
- 2) A training methodology using behavior cloning with an optimization based VRP expert solution for the coverage problem, as well as an extension to exploration with partially observable states.
- 3) A graph neural network architecture which explicitly respects equivariance in the task structure in order to achieve zero shot generalization to large maps and

Supported by ARL Grant DCIST CRA W911NF-17-2-0181, NSF Grant CNS-1521617, ARO Grant W911NF-13-1-0350, ONR Grants N00014-20-1-2822 and ONR grant N00014-20-S-B001, and Qualcomm Research. The first author acknowledges support from the NSF Graduate Research Fellowship.

¹ Dept. of Electrical and Systems Eng., University of Pennsylvania, USA eig@seas.upenn.edu

²Dept. of Mechanical Eng. and Applied Mechanics, University of Pennsylvania, USA

large teams for which an expert demonstration is intractable.

II. MULTI-ROBOT ROUTING FOR THE COVERAGE PROBLEM

We consider the time-constrained coverage path planning problem in which a team of robots must maximize the visited region of interest within a given time limit [13]. The problem of searching over the set of all feasible trajectories is intractable, so we finely discretize the region of interest, motivated by the success of motion planning in lattices [19]. In practice, a coarse map of a region may be provided via satellite imagery. To generate the lattice representation, we initialize a grid of points with a spacing of 5 meters, and then remove all points that are within an obstacle in the mission environment pictured in Fig. 1, and then add connections between adjacent nodes in free space.

We denote the set of robots as \mathcal{R} , waypoints as \mathcal{W} , and unvisited waypoints of interest denoted as \mathcal{X} . The set of unvisited waypoints of interest is a subset of all waypoints, $\mathcal{X} \subseteq \mathcal{W}$. Next, we define the parameters that describe the map of the environment: p^j is the location of waypoint j , and \mathcal{N}_j is the set of neighboring waypoints to waypoint j . Then, we define the mission: $x_t^j \in \{0, 1\}$ is an indicator of whether the waypoint j is of interest at time t , with $x_t^j = 1$ if $j \in \mathcal{X}$. T denotes the time budget for the mission. The location of robot i at time t is q_t^i and $\mathbb{1}_{q_t^i=p^j}$ indicates whether robot i is currently at waypoint j . The multi-robot coverage problem can now be formulated as:

$$\begin{aligned} \max_{\{q_t^i\}_{i \in \mathcal{R}}} & \sum_{t=0}^T \sum_{j \in \mathcal{W}} \sum_{i \in \mathcal{R}} x_t^j \mathbb{1}_{q_t^i=p^j} \\ \text{s.t.} & x_t^j = x_0^j \prod_{i \in \mathcal{R}} \prod_{s=0}^{t-1} (1 - \mathbb{1}_{q_s^i=p^j}), \quad \forall j \in \mathcal{W}, t \leq T \\ & q_{t-1}^i = p^j, q_t^i = p^k \Rightarrow k \in \mathcal{N}_j, \quad \forall i \in \mathcal{R}, \forall j, k \in \mathcal{W}, t \leq T. \end{aligned} \quad (1)$$

Alternate formulations exist that can pose the coverage problem as a Vehicle Routing Problem (VRP) to be solved as a mixed-integer program [20]. For a given problem instance, we can use existing routing solvers such as [17] to generate open loop solutions. Other approaches to the coverage problem typically pre-compute the roles and trajectories of all team members for the duration of the mission [13].

In contrast, our goal is to develop a closed-loop controller that computes only the next action for each robot based on the current state of the system. This approach generalizes easily to systems with real dynamics in which agents may not instantaneously transition between desired waypoints, such as the team of simulated robots in Fig. 1. It also enables generalization to dynamic graphs, permitting us to extend the coverage problem to an exploration problem in which new waypoints are discovered online during execution. We seek to learn a closed loop controller, π , that maximizes the observed waypoints in expectation over the set of initial

states and maps:

$$\begin{aligned} \max_{\pi} & \mathbb{E}_{q_0^i, p^j, x_0^j} \sum_{t=0}^T \sum_{j \in \mathcal{W}} \sum_{i \in \mathcal{R}} x_t^j \mathbb{1}_{q_t^i=p^j} \\ \text{s.t.} & x_t^j = x_0^j \prod_{i \in \mathcal{R}} \prod_{s=0}^{t-1} (1 - \mathbb{1}_{q_s^i=p^j}), \quad \forall j \in \mathcal{W}, t \leq T \\ & q_{t-1}^i = p^j, q_t^i = p^k \Rightarrow k \in \mathcal{N}_j, \quad \forall i \in \mathcal{R}, \forall j, k \in \mathcal{W}, t \leq T \\ & \{q_t^i\}_{i \in \mathcal{R}} = \pi \left(\{q_{t-1}^i\}_{i \in \mathcal{R}}, \{x_{t-1}^j\}_{j \in \mathcal{W}}, \{p^j\}_{j \in \mathcal{W}} \right). \end{aligned} \quad (2)$$

III. METHODS

To solve the problem of multi-robot coverage, we develop a parametrization of the problem as a heterogeneous graph that can be input to a Graph Neural Network to be trained via supervised learning.

A. Graph Representations for Coverage

We develop a parametrization of the current system state in which all entities, including robots and waypoints, are part of a single computation graph. Each robot in the team is considered as a node in the graph and robots are allowed to move from waypoint to waypoint in discrete steps. Due to robots' movement, the graph topology changes during execution. Using the lattice representation, we can abstract away the global positions of the robots, and maintain only relative information of nearby waypoints and robots. Unlike [11], we do not need to provide the relative distances of all waypoints to all agents. Furthermore, the action space of each robot is now discrete: the robot chooses one nearby waypoint to move to.

There are two types of edges: 1) map edges between waypoints to indicate free space and 2) action edges between robots and waypoints that indicate a robot's capability to move to nearby locations. Both types of edges are undirected. The waypoint connectivity is defined by the lattice induced over the given map of obstacles and free space. A robot is connected to the same waypoints as the waypoint at its current location, if $q_i = p_j$, then $\mathcal{N}_i = \mathcal{N}_j$.

The feature vector for each node \mathbf{v}_i of index i indicates the type of this node:

$$\mathbf{v}_i = \{\mathbb{1}_{i \in \mathcal{R}}, \mathbb{1}_{i \in \mathcal{W}}, \mathbb{1}_{i \in \mathcal{X}}\}, \quad (3)$$

where $\mathbb{1}$ is an indicator function.

We define \mathbf{e}_k as an edge feature vector for the directed edge of index k , with a sender node s_k and a receiver node r_k . For this task, we define \mathbf{e}_k to be the distance between the positions of nodes s_k and r_k :

$$\mathbf{e}_k = \|p_{s_k} - p_{r_k}\|. \quad (4)$$

The set of all edge features is $E = \{\mathbf{e}_k\}$, and the set of vertex features is $V = \{\mathbf{v}_i\}$, and we denote the graph that represents the state of the entire system as $\mathcal{G} = \{E, V\}$. At a given time t , the state of the multi-robot task is described by \mathcal{G}_t .

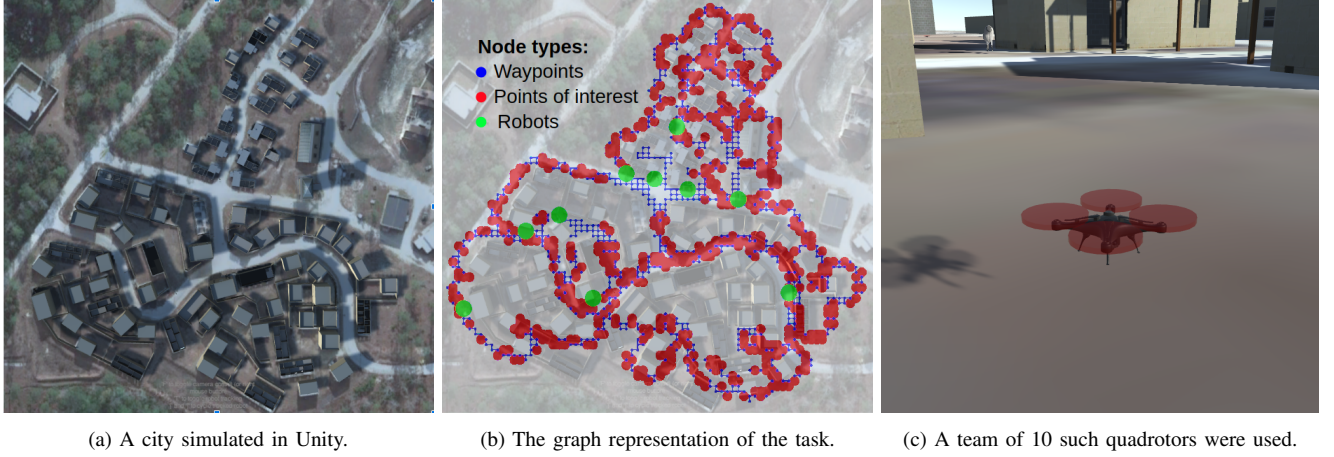


Fig. 1. The trained models were tested on a team of robots simulated in Unity and controlled by waypoint commands issued through a Robot Operating System interface. The trained model allows the robots to divide and conquer to visit the points of interest more efficiently than a greedy model. We visualize this experiment in a provided along with this work: <https://youtu.be/MiYSeNTyoA>

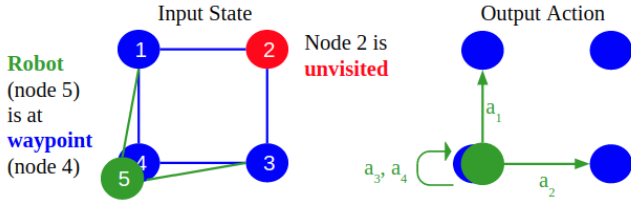


Fig. 2. Robots and waypoints comprise the nodes in the graph, with the edges between them indicating the ability of robots to move to new locations.

B. Graph Representations for Exploration

We view the exploration problem as the problem of coverage on a growing graph. Waypoint nodes are added to the graph when they are observed by a range sensor with range S : if $\|p_i^j - q_i^j\| \leq S$, then $\mathcal{W}_{i+1} = \mathcal{W}_i \cup \{p_i\}$, with the set of waypoints growing over time. Exploration introduces the possibility that an observed waypoint may or may not have adjacent waypoints that are currently unexplored. We call these frontier nodes and add an indicator feature to indicate whether a waypoint is part of the set of frontier nodes, \mathcal{F} :

$$\mathbf{v}_i = [\mathbb{1}_{i \in \mathcal{R}}, \mathbb{1}_{i \in \mathcal{W}}, \mathbb{1}_{i \in \mathcal{X}}, \mathbb{1}_{i \in \mathcal{F}}]. \quad (5)$$

C. Aggregation Graph Neural Networks

Graph Neural Networks are an increasingly popular tool for exploiting the known structure of any relational system [9]. In graph convolutional networks, the graph convolution operation is defined using learnable coefficients that multiply powers of the adjacency matrix times the graph signal [21], [22]. We extend this architecture by incorporating non-linear graph convolution operations.

The building block of a GNN is the Graph Network Block. Given a graph signal, $\mathcal{G} = \{\{\mathbf{e}_k\}, \{\mathbf{v}_i\}\}$, one application of the GN block transforms these features, $\mathcal{G}' = \{\{\mathbf{e}'_k\}, \{\mathbf{v}'_i\}\}$:

$$\mathbf{e}'_k = \phi^e(\mathbf{e}_k, \mathbf{v}_{r_k}, \mathbf{v}_{s_k}), \quad \mathbf{v}'_i = \phi^v(\bar{\mathbf{e}}'_i, \mathbf{v}_i), \quad \bar{\mathbf{e}}'_i = \rho^{e \rightarrow v}(E'_i). \quad (6)$$

$GN(\cdot)$ is a function of the graph signal \mathcal{G} , described by the application of ϕ^e , $\rho^{e \rightarrow v}$ and ϕ^v in that order to produce the transformed graph signal \mathcal{G}' , with the same connectivity but new features on the edges and nodes.

The aggregation operation $\rho^{e \rightarrow v}$ takes the set of transformed incident edge features $E'_i = \{\mathbf{e}'_k\}_{r_k=i}$ at node i and generates the fixed-size latent vector $\bar{\mathbf{e}}'_i$. Aggregations must satisfy a permutation invariance property since there is no fundamental ordering of edges in a graph. Also, this function must be able to handle graphs of varying degree, so the mean aggregation is particularly suitable to normalizing the output by the number of input edges [14]:

$$\rho^{e \rightarrow v}(E'_i) := \frac{1}{|E'_i|} \sum_{\mathbf{e}'_k \in E'_i} \mathbf{e}'_k. \quad (7)$$

The mean aggregation operation is especially helpful for improving the stability of GNNs with large receptive fields.

Next, we describe two variants of the Aggregation GNN architecture that build upon [23]. The linear Aggregation GNN architecture uses the following parametrization:

$$\phi_L^e(\mathbf{e}_k, \mathbf{v}_{r_k}, \mathbf{v}_{s_k}) := \mathbf{v}_{s_k}, \quad \phi_L^v(\bar{\mathbf{e}}'_i, \mathbf{v}_i) := \bar{\mathbf{e}}'_i, \quad (8)$$

while the non-linear Aggregation GNN uses learnable non-linear functions to update node and edge features:

$$\begin{aligned} \phi_N^e(\mathbf{e}_k, \mathbf{v}_{r_k}, \mathbf{v}_{s_k}) &:= \text{NN}_e([\mathbf{e}_k, \mathbf{v}_{r_k}, \mathbf{v}_{s_k}]), \\ \phi_N^v(\bar{\mathbf{e}}'_i, \mathbf{v}_i) &:= \text{NN}_v([\bar{\mathbf{e}}'_i, \mathbf{v}_i]), \end{aligned} \quad (9)$$

where NN_e and NN_v are 3 layer MLPs with 16 hidden units. Note that the linear Aggregation GNN in (8) cannot use the input edge features, such as those defined in (4), unlike the non-linear GNN defined in (9).

D. Policy Architecture

While a Graph Network Block can be used to compose a variety of architectures, for this work, we develop a variant of the Aggregation GNN in which the output of every GN

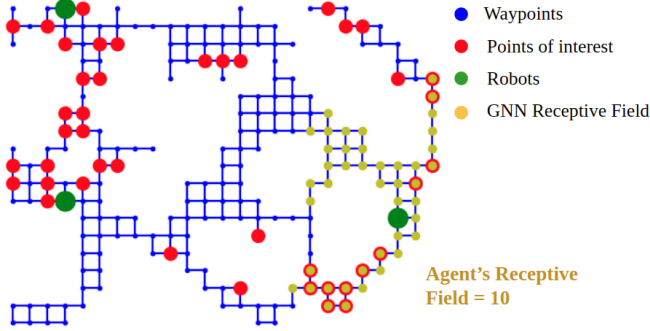


Fig. 3. The receptive field ($K=10$) of a Graph Neural Network.

stage is concatenated, and finally processed by a linear output transform [23]:

$$\mathcal{G}' = f_{\text{out}}\left(\left[f_{\text{dec}}(f_{\text{enc}}(\mathcal{G})), f_{\text{dec}}(GN(f_{\text{enc}}(\mathcal{G}))), \right. \right. \\ \left. \left. f_{\text{dec}}(GN(GN(f_{\text{enc}}(\mathcal{G}))), \dots \right]\right). \quad (10)$$

The addition of the encoder f_{enc} and decoder f_{dec} layers was inspired by the Encode-Process-Decode architecture presented in [9]. The number of GN operations is a hyperparameter and determines the *receptive field* (K), of the architecture, or how far information can travel along edges in the graph. A GNN with a receptive field of zero can be compared to the Deep Set architecture that neglects any relational data [24].

In our policy architecture, f_{enc} , f_{dec} are 3 layer MLPs with 16 hidden units, and a ReLU activation only after the first two layers. f_{out} is a linear function that reduces the high-dimensional latent space vectors output on the edges to the logits of a Boltzmann distribution. We sample from the edges that connect each robot node to neighboring waypoint nodes to determine the edge that each robot will next travel. This is in contrast to [14] where there was a hand-engineered feature extractor.

Fig. 3 visualizes a segment of a typical training scenario. Three robots shown in green must visit red points of interest by traveling along the waypoints in free space indicated by blue nodes and edges. We visualize the receptive field of a GNN with $K = 10$ indicating the 10-hop information available to the robot on the right. A larger receptive field allows each agent to use information about more distant regions of the map to compute the controller.

E. Baseline Controllers

The learned policies are compared to three types of controllers: 1) an expert open loop VRP solution, 2) a receding horizon VRP-based controller, and 3) a greedy controller. We use Google’s OR-Tools library [17] to provide optimization-based expert solutions to the VRP. This expert plans once for the full mission length of T and then this trajectory is executed in an open loop fashion. The expert assumes global knowledge of the map and may be intractable in larger or dynamic graphs. The training data was generated

with this open-loop approach. We also devise a receding horizon controller that plans for trajectories of $\hat{T} < T$, and then executes the first step of this trajectory, and then re-plans at the next time step. The Expert controller baselines use receding horizon control in Figs. 4, 5, 9. Finally, a greedy controller that routes each robot to the nearest unvisited point of interest is a great heuristic in many scenarios, so we include it as a practical lower bound. The greedy controller can be implemented with a finite receptive field using only a K -step distance matrix between nodes in the graph and we provide this benchmark in Figs. 4, 5, 9. A limited-horizon greedy controller may be more practical in larger graphs for which computing a full distance matrix is expensive.

F. Imitation Learning from Expert Demonstrations

In behavior cloning, our goal is to use stochastic gradient descent to minimize the difference between the expert’s action and the policy’s output, where \mathcal{L} is a cross-entropy loss, since the action space is discrete:

$$\pi^* = \underset{\pi}{\operatorname{argmin}} \sum_{(G_t, \mathbf{u}_t) \in \mathcal{D}} \mathcal{L}(\pi(G_t), \mathbf{u}_t). \quad (11)$$

To use behavior cloning to train the graph neural network policy, we require a dataset of expert trajectories. We collect a dataset of 2000 expert trajectories of length $T = 50$ in randomly generated graphs, $\mathcal{D} = \{(G_t, \mathbf{u}_t)\}_{t=1, \dots, 50}$. The graphs are generated by sampling regions of 228 waypoints on average from the graph shown in 1. The performance of the learned controller is tested on graphs generated from the same distribution over trajectories of length $T = 50$. The models were trained for 200 epochs with a batch size of 32. Adam optimizer was used with an initial learning rate of 0.001 which was decayed by a factor of 0.95 for every 200 batches.

For the exploration task, we use the expert controller that uses the full graph to generate the trajectory, but only the local state observations are stored in the dataset, so the robot is learning to predict what the omniscient centralized controller would do based only on partial observations of currently explored nodes and frontiers.

G. Implementation Details

We use DeepMind’s Graph Nets library and a variant of the Encode-Process-Decode architecture for the graph neural network policy [9].

We implement a local collision avoidance strategy as part of the task specification. A robot is allowed to move to a new waypoint if no other robot will be occupying this waypoint during the next time step. Conflicts are resolved by giving priority to the robot with the smaller index. If agents moved to the same waypoint, they are likely to continue to travel together since the policy cannot disambiguate the agents due to the graph parametrization. To avoid this redundancy, we eliminate the ability of robots to move to the same waypoint. In a real decentralized multi-robot team, collision avoidance based on on-board sensing may also be necessary.

Due to the known maximum degree of the map graph, we can fix the number of neighboring waypoints considered by the policy to be up to 4, so that existing infrastructure for stochastic policies with fixed-size action spaces can be used [25]. The output of the learned policy are weights for up to 4 edges from waypoints to robots, as shown in Fig. 2.

Finally, an open-source implementation of the learning architecture can be found here: https://github.com/katetolstaya/graph_rl and the applications here: <https://github.com/katetolstaya/gym-flock>.

IV. RESULTS

First, we highlight the impact of the Aggregation GNN’s receptive field on its performance on the exploration and coverage tasks. Next, we examine how the GNN can generalize to larger graphs and team sizes. Finally, we validate the use of the GNN controllers in a high-fidelity simulator.

A. Locality

On the coverage task, the learned controllers reliably outperform a greedy controller, but fall short of the receding-horizon expert in Fig. 4. For the GNNs, we see a sharp improvement in performance with increasing receptive field. The linear and non-linear variants of the aggregation GNNs perform comparably, with the non-linear GNN performing slightly better. With an increasing horizon, the greedy controller rapidly reaches its asymptotic mean reward of 70.1 with a standard error of the mean (SEM) of 1.32. The open-loop expert obtains a mean reward of 91.0 with a SEM of 0.87. In this experiment, the expert controller is optimal and obtains the maximum rewards possible and provides the upper bound on the performance of the GNN.

On the exploration task, the learned controllers outperform greedy and the planner-based controllers as we can see in Fig. 5. Existing planning-based solutions rely on ad-hoc heuristics to weigh the importance of a waypoint at the frontier versus other waypoints of interest, and we show that a learned solution can improve over a receding horizon expert controller. We also see a significant improvement in the mean reward obtained by the GNNs as their receptive field increases. The non-linear GNN is again slightly better than the linear GNN.

We further analyze the effect of a model’s receptive field on its performance in varying size graphs. A finite receptive field decomposes the problem into local neighborhoods for each robot, producing a decentralized solution. The robot only uses a fixed-hop neighborhood of the graph for computing the action as seen in Fig. 3. The graph diameter is the max distance between any two nodes in a graph. A model with a smaller receptive field performs worse than a model with a larger receptive field, especially in graphs of larger diameter as demonstrated by Fig. 8. As the graph size increases, there are more points of interest that must be visited. The larger receptive field controller is able to route the agents to these waypoints, while the controller with the

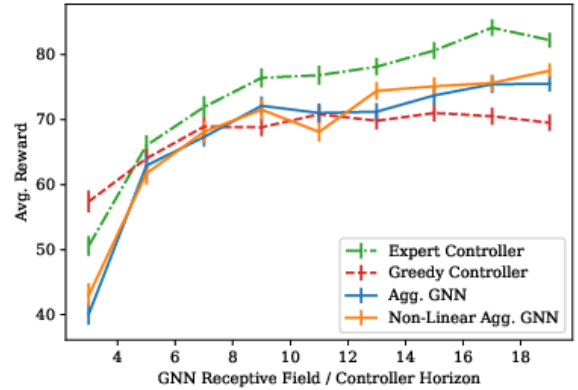


Fig. 4. GNNs with larger receptive fields achieve higher rewards on the coverage task. Mean reward over 100 episodes with standard error is shown.

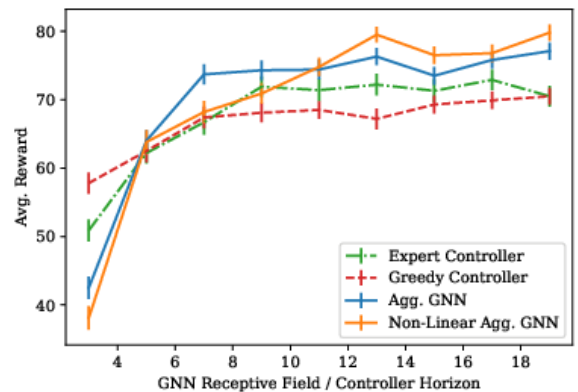


Fig. 5. GNNs surpass the expert controller on the exploration task. Mean reward over 100 episodes with standard error is shown.

smaller receptive field ($K = 3$) is unable to compute a high-reward path.

While the expert controller outperforms the learned controllers in Fig. 4, the VRP-based controller requires two orders of magnitude more time to compute a control action, as seen in Table I. The neural networks used a GTX 1080 GPU and the benchmarks used one core of an Intel Core i9-9900K processor. The expert and greedy controllers with receptive field of ∞ use the full adjacency matrix and, in particular, the expert plans one open-loop trajectory of length $T = 50$. We observe that the linear GNN has a faster controller computation time, and controllers with larger receptive fields also need more time.

B. Transference

The trained GNN models effectively generalize to larger robot team and map sizes than can be solved by conventional VRP solutions. The models were first trained on 4 agents and 228 waypoints on average. Then, the models were tested on a map size of 5659 waypoints with a graph diameter of 205. The team size varied from 10 to 100 agents. For both

Policy	Receptive Field		
	K=9	K=19	∞
Expert	13500	23500	2330
GN-MLP	176	277	-
GN-Linear	133	171	-
Greedy	86.3	142	297

TABLE I: The average controller computation time per episode in milliseconds, tested over 100 episodes of the coverage task. The expert controller (Receptive field = ∞) plans one open-loop trajectory for the task given the full adjacency matrix.

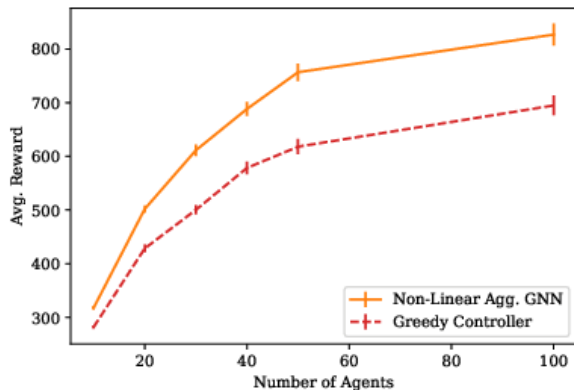


Fig. 6. Generalization to a coverage task with 5659 waypoints. We plot the average reward over 100 episodes with standard error. The GNN was trained with 4 agents and tested on teams of up to 100.

the coverage and exploration generalization experiments, the map and team sizes made the centralized expert solution intractable. In Fig. 6, we observe that the learned solution consistently outperforms the greedy controller on the coverage task. In Fig. 7, this difference is even bigger for the exploration task. We hypothesize that this is because the learned policy learns to weigh the frontier nodes more than other unexplored nodes.

C. Reinforcement Learning

Reinforcement Learning (RL) is also an effective tool for training Aggregation GNN controllers and achieves comparable performance to imitation learning in Fig. 9. RL is especially important in tasks that do not have suitable expert controllers for generation of training data. We use the Proximal Policy Optimization algorithm [26] and parametrize both the policy and value functions as Aggregation GNNs. The policy architecture is the same as for the imitation learning experiments, while the value function sums over the values output for the robot nodes. The GNN captures the known structure of the problem for the value and policy functions. The reward is the summation in the objective of (2), with no reward shaping required. The model was trained using 1×10^6 observations using an open source implementation of PPO from [25] using Adam optimizer with a step size

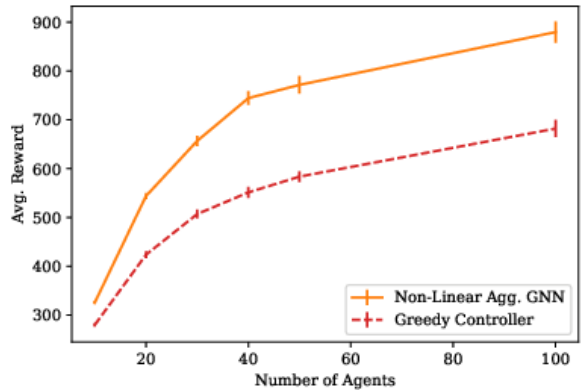


Fig. 7. Generalization to an exploration task with 5659 waypoints. We plot the average reward over 100 episodes with standard error. The GNN was tested on teams of up to 100.

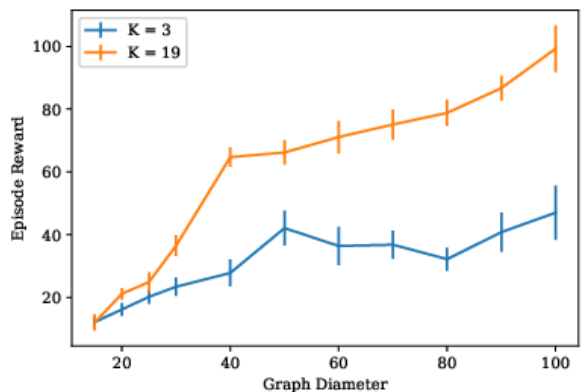


Fig. 8. Effect of receptive field of non-linear GNNs in graphs of varying diameters, as measured by the mean reward over 20 episodes and standard error.

of 1×10^{-4} and a batch size of 40. The controller trained using RL reached the performance of the imitation-trained controller for higher receptive fields. We observed a higher variance in the performance of the controllers trained with RL, so additional tuning of this system may be necessary.

D. Dynamics

Despite training on small-scale teams with instantaneous discrete state transitions, the GNN is effective for control in a coverage mission with ten quadrotors in a large simulated environment, pictured in Fig. 1. In 400 seconds of mission time, the team of 10 robots visited 490 points of interest using the greedy controller, as compared to visiting 610 points of interest using the non-linear GNN with a receptive field of 19. The use of the lattice representation with discrete states allows the model trained on an ideal discrete task to generalize zero-shot to a high-fidelity simulator. We discretize the robot positions provided by the simulator by clipping them to each robot’s nearest waypoint. The spatial aggregation operations are invariant to the time-scale of the

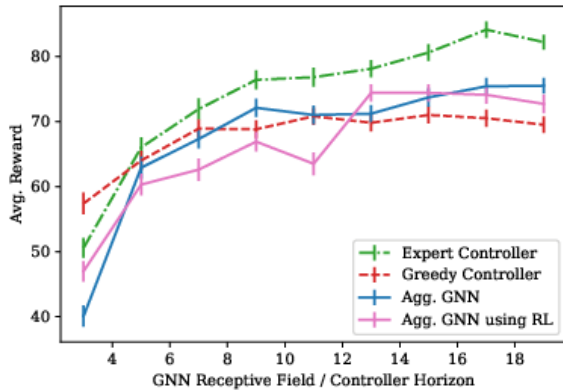


Fig. 9. Comparison of linear GNNs trained with imitation learning and reinforcement learning. Mean reward over 100 episodes with standard error.

task, so we can execute the GNN controller with the non-ideal dynamics of a quadrotor. The local collision avoidance strategy described in Section III-G was effective for this scenario because the implementation of the control policy and the evaluation of the GNN was centralized. Furthermore, each robot was assigned a different altitude. In more complex environments or for tasks specified on a 3D lattice, a policy that uses on-board sensing may be required.

V. CONCLUDING REMARKS

We develop a scalable GNN architecture for multi-robot coverage and exploration tasks. The approach surpasses existing decentralized heuristics and also scales well to from team sizes of 4 to teams of up to 100 agents. We also demonstrate that this architecture can be trained via reinforcement learning. As a bridge to deploying this approach to physical robot teams, we demonstrate generalization to a simulated robot team subject to dynamics in a dense urban environment.

To deploy the GNN in a real distributed team, we would need to address challenges such as asynchronous or intermittent communication. One possible solution could be the evaluation of the GNN on the contents of a robot's local buffer containing the estimated state of the system, and allowing intermittent communication among robots to update each other about the current positions of other robots, points of interest, and, for the exploration task, the growing map of waypoints. One approach to enable data distribution in mobile robot teams was explored in [27].

REFERENCES

- [1] V. Sharma, M. Bennis, and R. Kumar, "UAV-assisted heterogeneous networks for capacity enhancement," *IEEE Communications Letters*, vol. 20, no. 6, pp. 1207–1210, 2016.
- [2] S. Thrun, W. Burgard, and D. Fox, "A real-time algorithm for mobile robot mapping with applications to multi-robot and 3D mapping," in *Robotics and Automation, 2000. Proceedings. ICRA 2000. IEEE International Conference on*, vol. 1. IEEE, 2000, pp. 321–328.
- [3] S. Thrun and Y. Liu, "Multi-robot slam with sparse extended information filters," in *Robotics Research. The Eleventh International Symposium*. Springer, 2005, pp. 254–266.
- [4] B. Schlotfeldt, D. Thakur, N. Atanasov, V. Kumar, and G. J. Pappas, "Anytime planning for decentralized multirobot active information gathering," *IEEE Robotics and Automation Letters*, vol. 3, no. 2, pp. 1025–1032, 2018.
- [5] J. L. Baxter, E. Burke, J. M. Garibaldi, and M. Norman, "Multi-robot search and rescue: A potential field based approach," in *Autonomous robots and agents*. Springer, 2007, pp. 9–16.
- [6] J. S. Jennings, G. Whelan, and W. F. Evans, "Cooperative search and rescue with a team of mobile robots," in *Advanced Robotics, 1997. ICAR '97. Proceedings., 8th International Conference on*. IEEE, 1997, pp. 193–200.
- [7] H. Zhang and J. C. Hou, "Maintaining sensing coverage and connectivity in large sensor networks," *Ad Hoc & Sensor Wireless Networks*, vol. 1, no. 1-2, pp. 89–124, 2005.
- [8] B. Chen, B. Dai, and L. Song, "Learning to plan via neural exploration-exploitation trees," *arXiv preprint arXiv:1903.00070*, 2019.
- [9] P. W. Battaglia, J. B. Hamrick, V. Bapst, A. Sanchez-Gonzalez, V. Zambaldi, M. Malinowski, A. Tacchetti, D. Raposo, A. Santoro, R. Faulkner *et al.*, "Relational inductive biases, deep learning, and graph networks," *arXiv preprint arXiv:1806.01261*, 2018.
- [10] C. K. Joshi, T. Laurent, and X. Bresson, "An efficient graph convolutional network technique for the travelling salesman problem," *arXiv preprint arXiv:1906.01227*, 2019.
- [11] F. Chen, S. Bai, T. Shan, and B. Englot, "Self-learning exploration and mapping for mobile robots via deep reinforcement learning," in *AIAA Scitech 2019 Forum*, 2019, p. 0396.
- [12] J. Paulos, S. W. Chen, D. Shishika, and K. Vijay, "Decentralization of multiagent policies by learning what to communicate," in *2019 IEEE International Conference on Robotics and Automation (ICRA)*, Montreal, May 2019.
- [13] 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.
- [14] E. Tolstaya, F. Gama, J. Paulos, G. Pappas, V. Kumar, and A. Ribeiro, "Learning decentralized controllers for robot swarms with graph neural networks," in *Conference on Robot Learning*, 2020, pp. 671–682.
- [15] A. Khan, E. Tolstaya, A. Ribeiro, and V. Kumar, "Graph policy gradients for large scale robot control," in *Conference on Robot Learning*. PMLR, 2020, pp. 823–834.
- [16] P. Toth and D. Vigo, *The vehicle routing problem*. SIAM, 2002.
- [17] L. Perron and V. Furnon, "Or-tools," Google. [Online]. Available: <https://developers.google.com/optimization/>
- [18] Q. Sykora, M. Ren, and R. Urtasun, "Multi-agent routing value iteration network," *Int. Conf. on Machine Learning (ICML)*, 2020.
- [19] M. McNaughton, C. Urmson, J. M. Dolan, and J.-W. Lee, "Motion planning for autonomous driving with a conformal spatiotemporal lattice," in *2011 IEEE International Conference on Robotics and Automation*. IEEE, 2011, pp. 4889–4895.
- [20] G. Dantzig, R. Fulkerson, and S. Johnson, "Solution of a large-scale traveling-salesman problem," *Journal of the operations research society of America*, vol. 2, no. 4, pp. 393–410, 1954.
- [21] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *5th Int. Conf. Learning Representations*. Toulon, France: Assoc. Comput. Linguistics, 24–26 Apr. 2017.
- [22] F. Gama, A. G. Marques, G. Leus, and A. Ribeiro, "Convolutional neural network architectures for signals supported on graphs," *IEEE Trans. Signal Process.*, vol. 67, no. 4, pp. 1034–1049, Feb. 2019.
- [23] F. Gama, A. G. Marques, A. Ribeiro, and G. Leus, "Aggregation graph neural networks," in *44th IEEE Int. Conf. Acoust., Speech and Signal Process.* Brighton, UK: IEEE, 12–17 May 2019.
- [24] M. Zaheer, S. Kottur, S. Ravanbakhsh, B. Poczos, R. R. Salakhutdinov, and A. J. Smola, "Deep sets," in *Advances in Neural Information Processing Systems*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds., vol. 30. Curran Associates, Inc., 2017.
- [25] A. Hill, A. Raffin, M. Ernestus, A. Gleave, A. Kanervisto, R. Traore, P. Dhariwal, C. Hesse, O. Klimov, A. Nichol, M. Plappert, A. Radford, J. Schulman, S. Sidor, and Y. Wu, "Stable baselines," <https://github.com/hill-a/stable-baselines>, 2018.
- [26] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.
- [27] E. Tolstaya, L. Butler, D. Mox, J. Paulos, V. Kumar, and A. Ribeiro, "Learning connectivity for data distribution in robot teams," *arXiv preprint arXiv:2103.05091*, 2021.