

Multi-agent Reinforcement Learning with Graph Q-Networks for Antenna Tuning

Maxime Bouton, Jaeseong Jeong, Jose Outes, Adriano Mendo, Alexandros Nikou

Abstract—Future generations of mobile networks are expected to contain more and more antennas with growing complexity and more parameters. Optimizing these parameters is necessary for ensuring the good performance of the network. The scale of mobile networks makes it challenging to optimize antenna parameters using manual intervention or hand-engineered strategies. Reinforcement learning is a promising technique to address this challenge but existing methods often use local optimizations to scale to large network deployments. We propose a new multi-agent reinforcement learning algorithm to optimize mobile network configurations globally. By using a value decomposition approach, our algorithm can be trained from a global reward function instead of relying on an ad-hoc decomposition of the network performance across the different cells. The algorithm uses a graph neural network architecture which generalizes to different network topologies and learns coordination behaviors. We empirically demonstrate the performance of the algorithm on an antenna tilt tuning problem and a joint tilt and power control problem in a simulated environment.

Index Terms—multi-agent reinforcement learning, mobile networks, 6G, graph neural networks

I. INTRODUCTION

Mobile networks can be composed of thousands of base station antennas, and each of them comprises many parameters to be configured. The configuration of these parameters such as tilt or power usually has a great impact on the overall network performance. With the growing complexity of networks in 5G and beyond, the problem of dynamically configuring those parameters becomes increasingly challenging [1]. In addition to being costly, human monitoring and intervention on the network are not scalable. Multiple parameters can be used to affect the same performance metric, and changes in the configuration of one base station are likely to influence neighboring base stations and degrade their performance. In addition, the optimal choice of parameters is highly dependent on environmental factors as well as the spatial distribution of the traffic and the mobility of the connected user equipment. There is a need to design algorithms that can automatically learn tuning strategies to optimize the network performance by adapting to changes in the environment while considering all the possible interactions across neighboring base stations and across parameters.

In contrast to hand-engineered or manual tuning strategies, reinforcement learning (RL) provides a flexible framework to learn a control strategy from data. Previous works have demonstrated its application to a variety of radio access network optimization use cases including remote electrical tilt control,

optimization of handovers in 5G or power control, including industrial solutions [2]–[4]. Although they are promising, existing RL methods often fail to capture the needed coordination across neighboring base stations. For scalability reasons, they often resort to limiting assumptions and consider the control of one entity independently of the others, or rely on specific feature engineering to incorporate neighbor information as input to the reinforcement learning agents. In addition, they focus on optimizing local Key Performance Indicators (KPIs) involving one base station and its closest neighbors instead of considering the network as a whole in the reward formulation. In this work, we derive a multi-agent reinforcement learning algorithm capable to optimize network performance globally and control many base stations in a coordinated manner without relying on local reward approximation and feature engineering.

Our approach extends state-of-the-art cooperative multi-agent reinforcement learning algorithms by proposing a novel off-policy algorithm for cooperative learning when a graph structure is available. Inspired by value decomposition methods [5], [6], we propose a new Q-network architecture, graph Q-network (GQN), that is particularly well suited to address mobile network optimization problems. GQN uses graph neural networks to generalize to different network topologies and share knowledge across neighboring base stations and different parameters to control. Our GQN training algorithm uses a reward signal characterizing the performance of the network in a global area as opposed to using ad-hoc credit assignment techniques, and thereby, the agents learn to coordinate to improve the global network performance. We demonstrate the performance of the proposed algorithm on two mobile network optimization problems: tilt control and joint tilt and power control. In the first problem, we show that GQN outperforms all the baselines and can generalize to different network topologies. In the second problem, we show the ability to control different antenna parameters in a scalable way by considering each parameter as an agent. A technical report with additional details and experiments in appendix can be found at [preprinturl](#).

II. RELATED WORK

Existing works applying reinforcement learning for antenna tuning consider only local information of a cell and its closest neighbors [7], [8]. Several methods have been proposed to do so. The first one consists in hand-engineering an input feature and reward function to accommodate for the effect one cell can have on its neighbors. The KPIs of the neighboring cells are aggregated and the RL agent tries to optimize a combination of its own KPIs and the neighbor’s KPIs [9].

The authors are with Ericsson. Email: {maxime.bouton, jaeseong.jeong, jose.outes, adriano.mendo, alexandros.nikou}@ericsson.com

Other techniques have proposed to use graph neural networks to process neighbor information. However, they also require an ad-hoc engineering of the reward and can only control one base station at a time [10]. Other algorithms attempting to address the global network optimization problem have been proposed in previous works using coordination graphs [11]. This solution also required a heuristic to handle credit assignment between base stations by splitting individual rewards across neighbors. The inference cost of the message passing algorithm is larger than in our proposed method as it requires storing a neural network for every connected base stations in the graph. Our proposed algorithm can train a model controlling multiple antennas from a single global reward signal. In addition, previous works focus on controlling one antenna parameter, in this work we demonstrate the ability to control multiple parameters simultaneously (tilt and power).

Multi-agent RL algorithms have been proposed for problems where multiple agents cooperate for a common goal. The closest algorithms to our methods are value decomposition networks (VDN) and QMIX [5], [6]. They rely on factorizing the problem to only learn one value function per agent contrary to using individual reward signals. They train these individual value functions using one global reward such that the sum of individual values (or a weighted sum in the case of QMIX) matches the global reward. We take inspiration from these factorization methods but add a graph neural network component. It can exploit the topology of the telecommunications network to learn a more efficient decomposition of the joint action value function. Algorithms in the literature have not been applied to network optimization problems and have rarely been scaled to more than dozens of agents. Our algorithm scales to hundred of agents and can cope with a varying number of agents both at training and deployment time.

III. BACKGROUND

A. Multi-agent Reinforcement Learning

The problem of network optimization can be formulated as a multi-agent cooperative reinforcement learning problem where each network entity is an agent. The problem is modeled as a multi-agent Markov decision process [12]. Formally, it is described by the tuple $(\mathcal{S}, \mathcal{A}, T, R, \gamma)$, where \mathcal{S} is a joint state space, \mathcal{A} a joint action space, T an unknown transition model, $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ a global reward function, and $\gamma \in [0, 1)$ a discount factor. A joint state $s \in \mathcal{S}$ is equal to (s_1, \dots, s_n) where s_i is the state of agent i and n is the number of agents. Similarly, a joint action $\mathbf{a} \in \mathcal{A}$ is equal to (a_1, \dots, a_n) where a_i is the action of agent i . The agents do not need to be homogeneous and could have different action spaces. Our goal is to find a joint policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$ that maximizes the discounted accumulated global reward over time. A standard approach in RL is to represent π by a value function $Q : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ such that $\pi(s) = \arg \max_{\mathbf{a}} Q(s, \mathbf{a})$. In single-agent cases, the Q-learning algorithm can be used to learn Q . For problems with continuous, or large state spaces, Q can

be represented by a neural network and learned using deep Q-learning [13]. It minimizes the Bellman error:

$$\min_{\theta} E_{(s, \mathbf{a}, r, s')} [(Q(s, \mathbf{a}; \theta) - (r + \gamma \max_{\mathbf{a}'} Q(s', \mathbf{a}'; \theta)))^2] \quad (1)$$

where (s, \mathbf{a}, r, s') is a multi-agent experience sample collected by interacting with the environment and θ are the weights parameterizing the value function. The reward is a global reward for the whole system.

Representing and estimating Q in multi-agent systems is hard because of the dimensionality of the joint state and action spaces increasing with the number of agents. In addition, solving the maximization problem to find $\pi(s)$ is also challenging due to the large size of the joint action space. A common approach is to approximate the value function by factorizing it into individual value functions as follows [5]: $Q(s, \mathbf{a}) \approx \sum_{i=1}^n Q_i(s_i, a_i)$ where Q_i is the individual value function associated to agent i .

B. Graph Neural Networks

Graph neural networks (GNNs) are a family of neural network architectures designed to process graph structured data. Generally, a GNN is a differentiable parametric function that takes as input a graph with node and edge attributes. Consider a graph \mathcal{G} with vertices \mathcal{V} and edges \mathcal{E} where each node $i \in \mathcal{V}$ is associated to an attribute vector $x_i \in \mathbb{R}^d$. A simple graph neural network layer processing the input graph \mathcal{G} can be described by the following equation [14]: $h_i = \sigma \left(W_1 x_i + W_2 \sum_{j \in \mathcal{N}(i)} x_j \right)$ where h_i is the latent features of node i , σ is an activation function, W_1 and W_2 are learnable weight matrices, and $\mathcal{N}(i)$ represents the set of neighbors of vertex i in the graph. Other types of GNNs have been considered with different ways of aggregating neighbor features and considering edge features such as using convolution or attention like operators [14]–[16].

IV. PROPOSED APPROACH

In this section we describe how to model network optimization problems using cooperative multi-agent reinforcement learning. We then introduce a new algorithm that can efficiently exploit the inherent graph structure of the network to automatically coordinate different agents even in a scenario where they control different types of base station parameters.

A. System Model

Mobile networks are composed of several base stations equipped with antennas responsible for serving users in certain areas. An area of coverage and its associated antenna in a radio access network is referred to as a cell. In 5G and 6G networks, the number of cells in a network is expected to increase drastically by combining macro cells responsible for covering a large area and small cells providing increased capacity in targeted zones. Each of the antennas broadcasting a signal to the cell can be tuned to maximize the signal quality which results in improved quality of experience for the users of that cell. In this work we focus on two different antenna tuning scenarios: tuning the remote electrical tilt, and jointly tuning the tilt and

the downlink transmission power. Our goal is to demonstrate that our proposed algorithm can handle environments with heterogeneous parameters to be controlled.

Several metrics can be considered to maximize the user quality of experience such as throughput, coverage or signal quality. In this work, we focus on maximizing the downlink signal quality of all the users in the network as measured by the reference signal to interference and noise ratio (SINR) in dB. The goal of our algorithm is to find a policy mapping the global state of the network to a joint cell configuration in order to maximize the SINR of the whole network.

The SINR, ρ_u of a user u in a cell c depends on the tilt angle α and the downlink power ϕ as follows $\rho_u(\phi, \alpha) = \frac{R_{c,u}(\phi_c, \alpha_c)}{\sum_{i=1, i \neq c}^{N_{\text{cells}}} R_{i,u}(\phi_i, \alpha_i) + \mu}$ where $R_{c,u}$ is the reference signal received power (RSRP). The relation between RSRP and antenna parameters is derived through antenna models defined in standards [17]. More detailed on the RSRP calculation can be found in the appendix C of our technical report.

We consider the global SINR of the network as the average of the SINR of all the users in the network: $\text{SINR}_G(\alpha, \phi, \mathbf{s}) = 1/N_{\text{users}} \cdot \sum_u \rho_u$. Considering a geometric mean would have also been possible to model fairness among the users. The global SINR is a function of the joint configuration of all the cells in the network, along with some other exogenous factor such as the environment and the traffic distribution, all abstracted in the variable \mathbf{s} . In contrast to the global SINR, we define a local SINR for a given cell c : $\text{SINR}_{L,c} = 1/N_c \cdot \sum_{u \in c} \rho_u$ which is the average SINR of all the users connected to that cell, where N_c is the number of users connected to cell c .

B. Global Network Optimization as a Multi-agent MDP

Maximizing the global network performance can be formulated as a multi-agent Markov decision process where each cell in the network is an agent. We assume that a graph modelling the relation between the different cells is available. Each node is a cell, and a vertex exists between two cells if they can influence each other. Automatic neighbor relations methods defined in the standards can be used to identify these edges [17]. In our simulation, we use a criterion based on the geographical location and the azimuth angle of each cell along with the interference ratio between the cells. Namely, if cells are close to each other and have high mutual interference they are connected in the graph. An example of such a graph is illustrated in Fig. 1.

Similarly to previous work [10], each agent i is assumed to observe the following quantities: the position of the antenna (x_i, y_i) , the direction of the antenna (x_i^Δ, y_i^Δ) , the 10th, 50th, and 90th percentile of the SINR of the connected users, the antenna tilt angle α_i , and the antenna maximum downlink power ϕ_i . The joint state space \mathcal{S} is represented by the Cartesian product of the individual state of each cell.

We now consider two different problems: a tilt control problem and a joint tilt and power control problem. Each problem has a different action space and a different reward function. The true objective of these antenna tuning problems is to maximize the global SINR of the network to improve

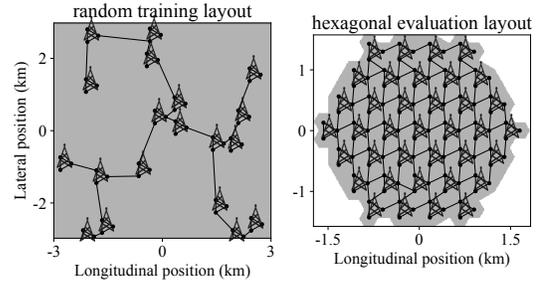


Fig. 1. Macro base station network deployments used in our experiments with the graph connecting the different cells (three per base station). The topology of the graph can change depending on the intersite distance. Evaluation layouts use denser deployments with more agents. The position of the node is offset from their real position for better visualization.

the quality of service and minimize the transmitted power to reduce the energy consumption (in the case where power is controllable). Existing multi-agent RL algorithms often must rely on ad-hoc decompositions of the reward signal in order to scale to a large number of agents. In order to compare with these baselines, we define both a global reward function that can be used by our method and a local reward function that will be used by the proposed baseline algorithm in our experimental section. The global SINR is defined as the arithmetic mean of the SINR of all the users in the network. However, other definitions such as geometric mean could be used in order to promote fairness. Our method could be applied in a straightforward way to these other definitions.

Tilt control: The tilt action space is defined by a set of remote electrical tilt changes of $\{-1^\circ, 0^\circ, +1^\circ\}$. The tilt is always bounded within the range $[0^\circ, 15^\circ]$. In this scenario, the maximum transmission power of each cell is fixed to 40 W. The local and global reward functions are defined as follows:

- Global SINR reward: $R(\mathbf{s}, \mathbf{a}, \mathbf{s}') = \text{SINR}_G(\alpha, \phi, \mathbf{s})$.
- Local SINR reward for cell i : $R_i(\mathbf{s}, \mathbf{a}, \mathbf{s}') = \text{SINR}_{L,i} + 1/N_i \cdot \sum_{j \in \mathcal{N}(i)} \text{SINR}_{L,j}$. This local formula considers the performance of a single cell and the average performance of its neighbors. We consider all the neighbors in the graph and each of them are equally weighted. A disadvantage of this approach is that we need to choose how to weigh the neighbor manually, while our proposed approach will learn automatically the reward decomposition. Similarly to the formula above, we only consider tilt when using this formula.

Joint tilt power control: The joint action space is defined by the Cartesian product of the tilt action space described above, and a set of maximum downlink power changes $\{-5 \text{ W}, 0 \text{ W}, +5 \text{ W}\}$ resulting in a total of 9 actions per cell. The power is bounded to the range $[10 \text{ W}, 60 \text{ W}]$. The reward functions are defined by:

- Global SINR and power reward $R(\mathbf{s}, \mathbf{a}, \mathbf{s}') = (1 - w) \text{SINR}_G(\alpha, \phi, \mathbf{s}) - w \cdot 1/N_{\text{cells}} \cdot \sum_c \phi_c$. This reward definition adds a penalty based on the average transmitted power of the cells. The penalty is weighted by a hyperparameter w . Setting w to 0 would yield to maximizing the power

(since it leads to maximum SINR) and optimizing tilt, setting w to 1 would lead to completely minimizing the power and disregarding any effect of the tilt. In practice, the choice of w should be driven by business intents from the network operators.

- Local SINR and power reward $R_i(\mathbf{s}, \mathbf{a}, \mathbf{s}') = (1 - w)\text{SINR}_{L,i} - w\phi_i + 1/N_i \cdot \sum_{j \in \mathcal{N}(i)} [(1-w)\text{SINR}_{L,j} - w\phi_j]$

Given a joint network state \mathbf{s} , we are interested in a policy giving a joint configuration of tilt and power to maximize the expected accumulated reward using the formalism of multi-agent reinforcement learning described in Section III-A. The number of actions per cell grows exponentially with the number of parameters considered. Since we consider only tilt and power, the size of the action space is still reasonably small and we can learn a single joint controller for tilt and power. In our experiments, we are also going to consider the possibility of making two separate agents for each cell, one controlling the tilt and one controlling the power.

Maximizing the global reward is generally challenging due to the complex interactions between the cells in the networks. Instead one can rely on value decomposition techniques such as independent Q-learning [18]. Each cell would be learning using their local observation and the proposed local reward definition considering only direct neighbors. The local reward ignores the fact that some neighbors might be more important than other and that indirect interactions with neighbors further away in the graph could affect performance. To consider the global effect of all the cell in the network, we propose the graph Q-network algorithm.

C. Graph Q-Network

We describe a multi-agent reinforcement learning algorithm to control multiple cells in the network optimizing a single global objective. The algorithm relies on graph neural networks to process information from all the agents. It then uses a factorization technique to learn individual value functions for each cell from a single global reward signal, rather than relying on hand-engineered reward decomposition across the cells. At deployment time, the trained model can be evaluated with an arbitrary number of cells allowing to train at small scale but deploy at large scale. These properties are mostly enabled by an original representation of the joint state action value function which we name graph Q-network (GQN).

To train such a model using a single global reward signal for all the agents, we rely on a factorization technique. Factorization generally consists of decomposing a very large function into a combination of smaller components. In this problem, we wish to learn an additive decomposition of the joint state action value function such that: $Q(\mathbf{s}, \mathbf{a}) = \sum_i Q_i(h_i, a_i)$ where $(h_1, \dots, h_n) = \text{GQN}(\mathbf{s}, A)$. The GQN function represents our proposed model and A denotes the adjacency matrix of the network graph. Intuitively, the GQN is going to learn a hidden state representation for each agent, such that the global value function is linearly factorizable into individual value functions in that hidden space. To simplify the notation, we note each

output node of the GQN as $Q_i(\mathbf{s}, A, a_i)$, which represents the value of taking action a_i for node i .

The GQN architecture is illustrated in Fig. 2. It is parameterized by weights and is end to end differentiable. The model takes as input a feature vector for each agent to control along with a graph representation of the mobile network. The construction of such graph is described in the previous section and has been also demonstrated in previous work [10], [11]. Each agent is represented by a node in the graph. The feature vector s_i corresponds to the local state of a cell as described in the previous section. It is processed by an encoding layer which consists of multilayer perceptrons (MLPs) applied individually to each node feature vector. To improve sample efficiency during training and the generalization of the model to an arbitrary number of nodes, the MLPs have shared weights.

The encoded features, h_i^e , are then processed by several GNN layers. We experimented with the GNN architectures from Morris *et al.* [14] and the graph attention neural networks from Veličković *et al.* [16]. The output of the GNN layers consists of another set of hidden features, h_i^d , one for each node, passed through a decoding layer. The model is trained such that the output node embeddings represent the individual value function of each agent which we note as: $Q_i(h_i^d, a_i)$.

Training a GQN model follows an off-policy training procedure similar to DQN and can benefit from existing innovations like prioritized experience replay, target networks and double Q learning [19]. The convergence guarantees of GQN are similar to the one provided by DQN. The training procedure of GQN is described in our technical report.

The GQN architecture has the sufficient property (additivity) to satisfy the individual global maximum principle [20] which states that the jointly optimal actions are given by taking the individual maximum of each individual value function at the output of the GQN model:

$$\mathbf{a}^* = \arg \max Q(\mathbf{s}, \mathbf{a}) = \begin{bmatrix} \arg \max Q_1(h_1^d, a_1) \\ \arg \max Q_2(h_2^d, a_2) \\ \vdots \\ \arg \max Q_n(h_n^d, a_n) \end{bmatrix} \quad (2)$$

Contrary to previous work [10], the model outputs state-action value functions for all the agents at the same time.

V. EXPERIMENTS

We empirically demonstrate the performance of our proposed algorithm in simulated environments. Using the scenarios in Section IV-B, we demonstrate three different experiments illustrating respectively: the training performance for tilt environment, the generalization performance for tilt control, the training performance for the joint control problem with an additional agent decomposition scheme.

A. Experiment Setup

Our experiments rely on a proprietary system level simulator which implements a ray tracing propagation model for computing the path gains at various user locations [21]. Details

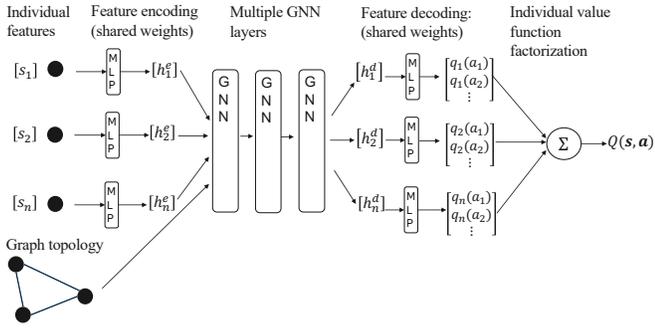


Fig. 2. Graph Q-network architecture

about the simulation parameters and training parameters can be found in our technical report in appendix¹.

We compare the following algorithms:

- DQN: The DQN algorithm with a target network, double Q-learning, prioritized replay, and distributed training [19].
- Neighbor DQN (N-DQN): it is a simple extension to DQN where the observation is augmented with the observation of the neighboring cells, proposed in previous work [10].
- GAQ [10]: the graph attention Q-network algorithm processes neighbor features using a graph attention layer. Contrary to our proposed approach, GAQ is trained using a local reward and is not able to perform joint control.
- GQN: our proposed method. We write GQN(GAT) when we used a graph attention layer [16], otherwise we used a graph convolutional layer from Morris *et al.* [14].
- Heuristic (H): This method is a rule-based method that sets the tilt angle such that the beam points to the middle of the cell. It observes the height of the antenna and the distance to its closest neighboring cell and calculates a desired tilt angle.

We did not add QMIX as part of the baselines as it is already shown to be outperformed by GAQ in previous work [10]. On the joint tilt and power control problem, QMIX was at most as good as DQN for some seeds, but the training was too unstable to report meaningful results in this paper. In the figures, the solid lines represent the mean over 3 random seeds and the shaded area is the 95th percentile confidence interval. The line for the heuristic represents the average performance and confidence interval evaluated on 300 episodes.

B. Training performance

Figure 3 illustrates the performance of GQN, GQN(GAT), DQN, GAQ and N-DQN on a tilt tuning scenario. We report the average SINR of the whole network (averaged over all the user equipments) which corresponds to the global SINR optimization target. The results are shown when training in a hexagonal deployment. A training step corresponds to one interaction with the environment. Similar conclusions can be drawn when training in the scenario with random deployments.

¹Experimental details can be found in appendix A here: technicalreport

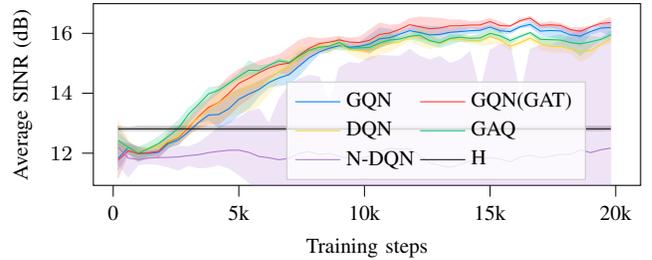


Fig. 3. Evolution of the average SINR per episode during the training of the agents for tilt control. We compare our approach (GQN, GQN(GAT)) to a standard DQN approach, N-DQN, and GAQ, a state-of-the-art algorithm developed for antenna tilt control [10].

We can see that all the methods except from N-DQN have similar convergence properties. In these experiments, all the methods follow DQN-like procedure and the convergence was mostly affected by the choice of the exploration schedule which is the same for all the baseline for a fair comparison. For the heuristic, we plot the average performance, as it is not a learning based algorithm and the performance is constant. We can see that all the RL algorithms outperform the baseline except for N-DQN.

While almost all the methods converge to an improved average SINR compared to the initial performance, GQN provides the best solution. We notice a difference of about 1 dB in the average performance. The convergence rate is roughly similar for all the methods with a slight disadvantage for GQN with the graph convolutional layer. The neighbor DQN approach performed poorly in general. We notice that the poor performance is accompanied by a high standard deviation. This is explained when looking to each of the three random seeds, only one performed at a level comparable to GAQ and DQN while the other seeds did not converge. We blame this instability on the complexity of the observation space and the difficulty to process this information with a simple fully connected neural network. The complexity of adding the graph neural network layers does not seem to affect the convergence for both GQN and GAQ. This result confirms the initial intuition that being able to train on the global objective rather than manually decomposing the global reward into individual cell reward (for DQN and GAQ) leads to better global performance.

C. Generalization to New Deployments

In this experiment, we save the trained models from Fig. 3 and evaluate them on network deployments unseen at training time. In addition, we evaluate another batch of models that were trained on random deployments rather than hexagonal deployments. We have two possible training configurations, hexagonal or random, and two evaluation configurations, hexagonal or random. Examples of a random training configuration and a hexagonal evaluation configuration are given in Fig. 1. We evaluate the models with up to 37 macro base stations which makes a total of 111 agents against 57 at training time.

Figure 4 illustrates the performance of the trained model on the hexagonal evaluation scenarios when training with random

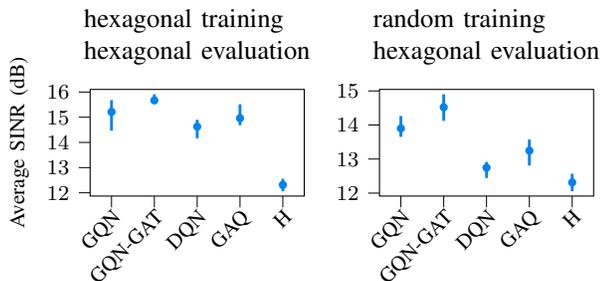


Fig. 4. Generalization performance of the GQN algorithm and the baselines on deployments unseen at training time. The points represent the average performance and the error bars represent the 95th percentile confidence interval (sometimes smaller than the marker).

or hexagonal deployments. The performance is average over 50 episodes and across three random seeds. We can see that the GQN methods provide overall better performance than the other methods, followed closely by the GAQ algorithm. Using a graph attention layer in GQN provides an additional gain in performance. An interesting outcome is that the GQN(GAT) model trained on random deployment, has a comparable performance on the hexagonal evaluation scenario (right plot) than the DQN and GAQ models that were trained specifically on hexagonal topologies (left plot). When the difference between training deployment and evaluation deployment is large (from random to hexagonal), the heuristic has a very close performance to some of the reinforcement learning algorithms, indicating that generalization to those scenarios is a difficult task. We omitted N-DQN from the figure for readability, its performance is around 9 dB in all cases. A poor performance is expected for this baseline as the neighbor features are stacked in a single vector. This input does not take into account the graph structure and is dependent on the order of the neighbors are stacked in.

D. Joint Control Problem

In this experiment, we consider the possibility to control heterogeneous agents. In the simplest case, the joint tilt and power control can be addressed by designing one agent per cell with an action space of size 9. This is the approach we will refer to as "joint". Such approaches might not be scalable if the number of parameters to control grows. Instead, one could consider two agents per cell: one controlling the tilt, and one controlling the power, each with an action space of size 3. This is the approach we use in the GQN algorithm in this experiment. Instead of controlling 57 agents, the GQN will control 114 agents. We simply modify the graph topology such that each cell consists of two nodes connected to each other, but with the same neighbors as in the joint control case.

Results are presented in Fig. 5. By looking at the two objectives, SINR and average power, we can see that the final SINR is similar for all the methods while the final power value shows a significant difference. For a similar average SINR, the GQN algorithm is able to reduce the power by more than 35%. In addition, the GQN methods tend to converge faster on this

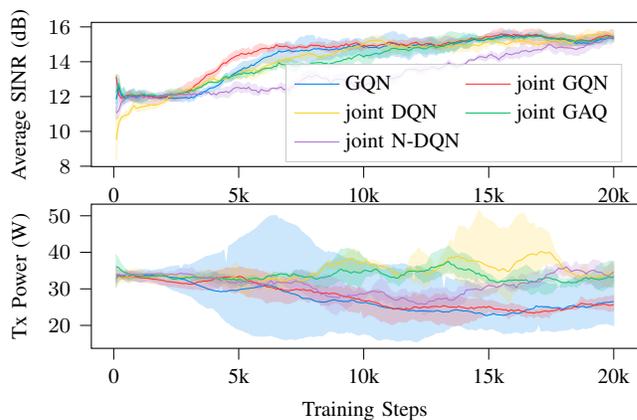


Fig. 5. Training performance of GQN with separate agents for tilt and power, joint GQN and the baselines in the joint control environment. For a similar average SINR, GQN is able to yield a greater power reduction.

problem. Our hypothesis is that the increase of dimensionality in the action space impacts more the methods relying only on fully connected neural networks.

Another observation is that the GQN approach yields the same average performance as joint GQN. This result suggests that GQN might be a good algorithm to decompose the antenna tuning problem into multiple agents per parameters. We expect that the standard deviation could be further reduced by more intense hyperparameter tuning.

VI. CONCLUSIONS

We presented a novel multi-agent reinforcement learning solution to tune many antennas in mobile networks. The proposed algorithm is an off-policy algorithm that, given a reward associated to the global network performance, learns to assign credits to the different antennas. Our method is able to learn the credit assignment and coordination behavior thanks to a graph Q-network, a graph neural network representation of the joint state action value function. Our experiments evaluate the learning algorithm in a tilt control scenario and a joint power and tilt control scenario. The results show that the trained model is able to learn control policies leading to a better global average SINR and more power savings. In addition, the trained GQN is able to generalize to denser network topologies unseen at training time, with almost double the number of agents. In the joint control scenario, we illustrated a way to separate the control of two parameters as different agents connected in a graph. By splitting the agents per parameter it allowed to maintain a smaller action space per agent while leading to similar performance as a joint controller.

Future works involve evaluating the algorithm on a broader range of problems and learning the graph structure automatically. We also consider investigating the control of advanced antenna systems which would include a much larger number of parameters per antenna and hence increase the number of agents in the graph.

REFERENCES

- [1] S. S. Mwanje, G. Decarreau, C. Mannweiler, M. N. ul Islam, and L. Schmelz, "Network management automation in 5g: Challenges and opportunities," in *IEEE International Symposium on Personal, Indoor and Mobile Radio Communications (PIMRC)*, 2016.
- [2] F. Vannella, J. Jeong, and A. Proutiere, "Off-policy learning in contextual bandits for remote electrical tilt optimization," *IEEE Transactions on Vehicular Technology*, 2022.
- [3] V. Yajnanarayana, H. Rydén, and L. Hévízi, "5g handover using reinforcement learning," in *IEEE 5G World Forum, 5GWF*, 2020.
- [4] E. Ghadimi, F. D. Calabrese, G. Peters, and P. Soldati, "A reinforcement learning approach to power control and rate adaptation in cellular networks," in *IEEE International Conference on Communications (ICC)*, 2017.
- [5] P. Sunehag, G. Lever, A. Gruslys, W. M. Czarnecki, V. F. Zambaldi, M. Jaderberg, M. Lanctot, N. Sonnerat, J. Z. Leibo, K. Tuyls, and T. Graepel, "Value-decomposition networks for cooperative multi-agent learning based on team reward," in *Autonomous Agents and Multiagent Systems*, 2018.
- [6] T. Rashid, M. Samvelyan, C. S. de Witt, G. Farquhar, J. N. Foerster, and S. Whiteson, "QMIX: monotonic value function factorisation for deep multi-agent reinforcement learning," in *International Conference on Machine Learning (ICML)*, 2018.
- [7] E. Balevi and J. G. Andrews, "Online antenna tuning in heterogeneous cellular networks with deep reinforcement learning," *IEEE Transactions on Cognitive Communications and Networking*, vol. 5, no. 4, pp. 1113–1124, 2019.
- [8] H. Farooq, A. Imran, and M. Jaber, "AI empowered smart user association in LTE relays hetnets," in *IEEE International Conference on Communications Workshops*, 2019.
- [9] F. Vannella, G. Iakovidis, E. A. Hakim, E. Aumayr, and S. Fegghi, "Remote electrical tilt optimization via safe reinforcement learning," in *IEEE Wireless Communications and Networking Conference*, 2021.
- [10] Y. Jin, F. Vannella, M. Bouton, J. Jeong, and E. A. Hakim, "A graph attention learning approach to antenna tilt optimization," in *International Conference on 6G Networking (6GNet)*, 2022.
- [11] M. Bouton, H. Farooq, J. Forgeat, S. Bothe, M. Shirazipour, and P. Karlsson, "Coordinated reinforcement learning for optimizing mobile networks," *NeurIPS Workshop on cooperative AI*, 2021. arXiv: 2109.15175.
- [12] M. J. Kochenderfer, *Decision making under uncertainty: Theory and application*. MIT Press, 2015.
- [13] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. A. Riedmiller, A. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [14] C. Morris, M. Ritzert, M. Fey, W. L. Hamilton, J. E. Lenssen, G. Rattan, and M. Grohe, "Weisfeiler and leman go neural: Higher-order graph neural networks," in *AAAI Conference on Artificial Intelligence*, 2019.
- [15] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *International Conference on Learning Representations (ICLR)*, 2017.
- [16] P. Velickovic, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio, "Graph attention networks," in *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*, 2018.
- [17] 3GPP, "Evolved Universal Terrestrial Radio Access (E-UTRA); Further advancements for E-UTRA physical layer aspects," 3rd Generation Partnership Project (3GPP), Technical Specification (TS) 36.814, 2017, Version 9.2.0.
- [18] M. Tan, "Multi-agent reinforcement learning: Independent versus cooperative agents," in *International Conference on Machine Learning (ICML)*, 1993.
- [19] M. Hessel, J. Modayil, H. van Hasselt, T. Schaul, G. Ostrovski, W. Dabney, D. Horgan, B. Piot, M. G. Azar, and D. Silver, "Rainbow: Combining improvements in deep reinforcement learning," in *AAAI Conference on Artificial Intelligence*, 2018.
- [20] K. Son, D. Kim, W. J. Kang, D. Hostallero, and Y. Yi, "QTRAN: learning to factorize with transformation for cooperative multi-agent reinforcement learning," in *International Conference on Machine Learning (ICML)*, 2019.
- [21] H. Asplund, M. Johansson, M. Lundevall, and N. Jaldén, "A set of propagation models for site-specific predictions," in *12th European Conference on Antennas and Propagation (EuCAP 2018)*, 2018.
- [22] E. Liang, R. Liaw, R. Nishihara, P. Moritz, R. Fox, K. Goldberg, J. Gonzalez, M. I. Jordan, and I. Stoica, "Rllib: Abstractions for distributed reinforcement learning," in *International Conference on Machine Learning (ICML)*, 2018.
- [23] M. Fey and J. E. Lenssen, "Fast graph representation learning with PyTorch Geometric," in *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2019.

APPENDIX

The appendix contains additional details about the method, parameters used for the experiments and additional experiments.

A. EXPERIMENTAL DETAILS

In this appendix, we describe the different hyperparameters for the algorithm and some implementation details about the experiments. Our simulator models an LTE network with antennas that can be controlled remotely for changing the electrical downtilt angle and the maximum downlink transmitted power and operating at a frequency of 2 GHz. The network has 10 000 users uniformly distributed on the map, generating an average traffic of 1 Mbps per cell.

The training consists of 20 000 steps split into episodes of 20 steps. At the beginning of an episode, a new deployment is sampled. We consider 19 base stations in the training environment, each consisting of 3 antennas, which makes a total of 57 cells. Each cell is associated to a learning agent. The average intersite distance between base stations is uniformly sampled between 300 m and 1500 m. We consider both hexagonal deployment and randomly generated deployment of the base stations with a minimum intersite distance, as illustrated in Fig. 6. At the beginning of an episode, the electrical tilt of each antenna is reset to a random value within the range $[0^\circ, 15^\circ]$. In the joint tilt and power control experiment the power is also reset to a random value while for the tilt experiment it is fixed at 40 W.

The baselines are using exactly the same simulator configuration and are trained with the same random seed such that they experience exactly the same simulated networks both in training and evaluation. The observation features and the reward signal are normalized during training. They all use an ϵ -greedy policy for exploration with a decay of the exploration rate, ϵ , from 1 to 0.01 during the first half of the training. Each training is repeated with three different random seeds.

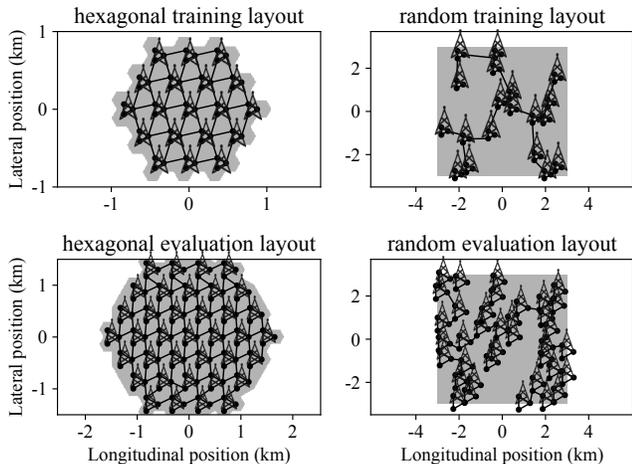


Fig. 6. Visualization of the hexagonal and random deployments used for training and evaluation. Macro base station network deployments used in our experiments with the graph connecting the different cells (three per base station). The topology of the graph can change depending on the intersite distance. Evaluation layouts use denser deployments with more agents. The position of the node is offset from their real position for better visualization.

TABLE I
HYPERPARAMETERS FOR TRAINING THE BASELINES AND OUR ALGORITHM.
THE FIRST LINES ARE COMMON TO ALL ALGORITHMS.

γ	0.0
learning rate	1×10^{-3}
initial ϵ	1.0
final ϵ	0.01
ϵ decrease steps	10000
activation function	ReLU
batch size	64
(N)DQN architecture	FC(64), FC(32)
GAQ architecture	GAT(32, 6 heads), GAT(32, 6 heads), FC(32), FC(32)
GQN architecture	GCN(32), GCN(32), FC(32), FC(32)
GQN(GAT) architecture	GAT(32, 4 heads), GAT(32, 4 heads), FC(32), FC(32)
GQN learning rate	1×10^{-2}

All the baselines are implemented using rllib [22] and the Pytorch geometric library [23], they all use a target network, prioritized replay, double Q-learning and distributed experience collection with 5 workers, each using 3 CPUs. We report all the hyperparameters in Table I. For all the baselines we carried out a hyperparameter search on the learning rate and number of layers and we rescaled the observation vectors and reward signals such that their value is between $[-1, 1]$. We used a discount factor of 0.0 since our problem’s goal, reaching the final optimal antenna configuration, does not require visiting intermediate states with performance degradation. Increasing the discount factor to 0.9 did not lead to any benefits for any of the baselines for this specific application. However, the method, as described in the previous section, still holds for larger discount factors. It leads to the same asymptotic results but with slower convergence.

In the N-DQN baseline, a maximum of 5 neighbors is

considered. Each observation vector from the neighbors is stacked and fed to a feed forward neural network.

For our GQN algorithm, we tried adding encoding MLPs before the GNN layer but it did not bring any improvement in performance (nor did it damage it). We also experimented with a graph convolutional neural network layer instead of the graph attention one.

We provide a pseudocode of our proposed GQN method in Algorithm 1. It follows a similar off policy training procedure as DQN. We use a target network, double Q learning and prioritized experience replay. The state of each agent and the graphs are stored in the replay buffer.

Algorithm 1 GQN training procedure.

- 1: **Initialize:** GQN weights, replay buffer $\mathcal{D} = \{\}$, training steps T , batch size B , ϵ schedule
 - 2: **for** $t = 1, \dots, T$
 - 3: Observe joint state \mathbf{s} , and adjacency matrix A
 - 4: With probability ϵ , choose a random joint action \mathbf{a}
 - 5: Otherwise choose $\mathbf{a} = \mathbf{a}^*$ according to Eq. (2).
 - 6: Observe the next joint state \mathbf{s} , the next adjacency matrix A' , and the global reward r .
 - 7: Store the transition $(\mathbf{s}, A, \mathbf{a}, r, \mathbf{s}', A')$ in \mathcal{D} .
 - 8: Sample batch $\{(\mathbf{s}^k, A^k, \mathbf{a}^k, r^k, \mathbf{s}'^k, A'^k)\}$ for $k = 1, \dots, B$ from \mathcal{D} .
 - 9: Assign $y^k = r^k + \gamma \max_{\mathbf{a}'} \sum_i Q_i(\mathbf{s}'^k, A'^k, a_i^k) \forall k$ as per Eq. (1).
 - 10: Perform a gradient descent step on the loss $1/B \cdot \sum_k [(\sum_i Q_i(\mathbf{s}^k, A^k, a_i^k) - y^k)^2]$
 - 11: **end for**
-

B. ADDITIONAL EXPERIMENTS

Tilt Generalization

In the tilt generalization experiment, we tried different combination of training the models on random deployments and evaluating on hexagonal deployments and vice versa. The results are presented in Fig. 7. In all cases, the GQN-GAT method present the best performance. An interesting outcome is that the GQN(GAT) model trained on random deployment, has a better performance on the hexagonal evaluation scenario (bottom left plot) than the DQN and GAQ models that were trained specifically on hexagonal topologies (top left plot). On the evaluation with random cell layout, the heuristic has a very close performance to some of the reinforcement learning algorithms, indicating that generalization to those scenarios is a difficult task, especially when the difference between training and deployment layout is large (top right plot).

Joint Tilt and Power Control

In this scenario the reward function is parameterized by w which controls the trade-off between minimizing power and maximizing signal quality. We experimented with value of w in $\{0.05, 0.1, 0.15, 0.2, 0.5\}$ for the reward functions. For the most extreme values for w the agents learns to set the power to the maximum value (low w), or the minimum (high w).

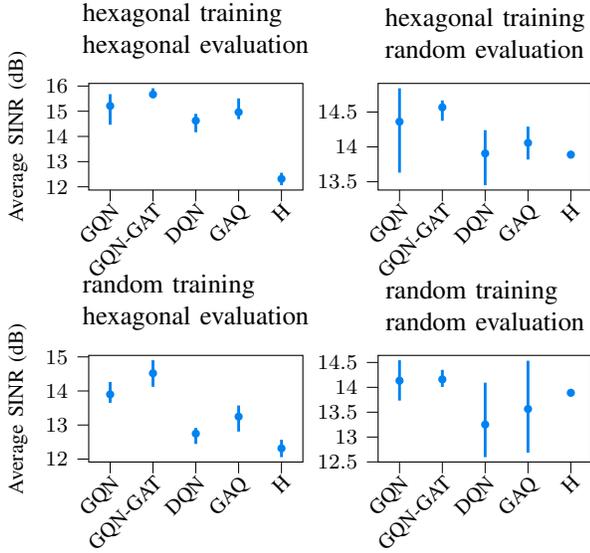


Fig. 7. Generalization performance of the GQN algorithm and the baselines on deployments unseen at training time. The evaluation deployments are denser, as illustrated in Fig. 1. The points represent the average performance and the error bars represent the 95th percentile confidence interval (sometimes smaller than the marker).

The values of 0.15 and 0.2 gave the most interesting results where the agent converged to non extreme power values. We witnessed the same effect for all methods and only report the results for $w = 0.15$ for the sake of readability.

In Fig. 8, we add the evolution of the global reward. The GQN model are directly trying to optimize this reward while other algorithms must rely on local signals. As a consequence they do not give as much power reduction as GQN.

C. RSRP AND SINR CALCULATIONS

To derive the relation between SINR and cell configurations we first define the reference signal received power (RSRP) $R_{c,u}$ for a user u connected to a cell c :

$$R_{c,u}(\phi, \alpha) = \frac{P_c(\phi)G_{c,u}(\alpha)}{L_{c,u}} \quad (3)$$

where $P_c(\phi)$ is the transmitted power of the antenna per reference signal resource element, as a function of the maximum transmitted power ϕ . $G_{c,u}(\alpha)$ is the antenna gain which depends on the tilt angle α and the azimuth which is kept fixed in our model. $L_{c,u}$ is the path loss. To compute the antenna gains, one can use the relationship defined in the third generation partnership project [17]. The maximum power ϕ

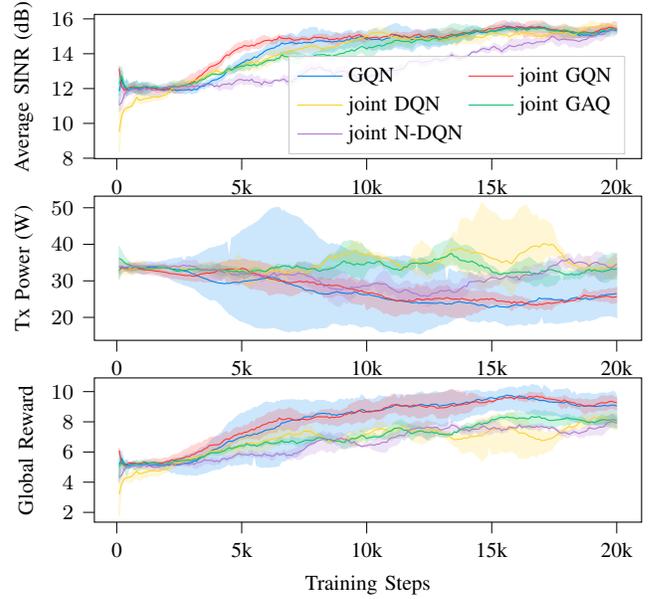


Fig. 8. Training performance of GQN (including global reward) with separate agents for tilt and power, joint GQN and the baselines in the joint control environment. For a similar average SINR, GQN is able to yield a greater power reduction. The solid lines represent the mean over 3 random seeds and the shaded area is the 95th percentile confidence interval.

is divided by the number of available resource blocks and multiplied by the cell specific reference signal power boost gain (set to one in our simulation). The received power is only depending on the cell that the user is connected to and on the propagation environment which affects the path loss. The user also receives power from other cells which is regarded as interference in the SINR calculation. For a user u , the cell c that the user is attached to is assumed to be the cell yielding the largest received power. The downlink SINR of a user u is defined as the ratio between the received power from cell c and the sum of the received power from all other cells, and the noise power μ :

$$\rho_u(\phi, \alpha) = \frac{R_{c_u}(\phi_c, \alpha_c)}{\sum_{i=1, i \neq c}^{N_{\text{cells}}} R_{i,u}(\phi_i, \alpha_i) + \mu} \quad (4)$$

The noise power is calculated over a frequency bandwidth of one resource element (15 kHz). Improving the SINR of a user involves improving the received power as well as reducing interference from other cells.