# Evaluating Generalization and Transfer Capacity of Multi-Agent Reinforcement Learning Across Variable Number of Agents

**Bengisu Guresti,**[1] **Nazim Kemal Ure** [1]

[1] Istanbul Technical University
guresti15@itu.edu.tr, ure@itu.edu.tr

## Abstract

Multi-agent Reinforcement Learning (MARL) problems often require cooperation among agents in order to solve a task. Centralization and decentralization are two approaches used for cooperation in MARL. While fully decentralized methods are prone to converge to suboptimal solutions due to partial observability and nonstationarity, the methods involving centralization suffer from scalability limitations and lazy agent problem. Centralized training decentralized execution paradigm brings out the best of these two approaches; however, centralized training still has an upper limit of scalability not only for acquired coordination performance but also for model size and training time. In this work, we adopt the centralized training with decentralized execution paradigm and investigate the generalization and transfer capacity of the trained models across variable number of agents. This capacity is assessed by training variable number of agents in a specific MARL problem and then performing greedy evaluations with variable number of agents for each training configuration. Thus, we analyze the evaluation performance for each combination of agent count for training versus evaluation. We perform experimental evaluations on predator prey and traffic junction environments and demonstrate that it is possible to obtain similar or higher evaluation performance by training with less agents. We conclude that optimal number of agents to perform training may differ from the target number of agents and argue that transfer across large number of agents can be a more efficient solution to scaling up than directly increasing number of agents during training.

## Introduction

A significant amount of problems in Multi-agent Reinforcement Learning (MARL) require agents to achieve a common goal, which makes learning to cooperate crucial. Nevertheless, as the number agents that need to cooperate increases, the difficulty of achieving optimal cooperation also increases. In order to create large-scale MARL applications, it is critical to find high performance solutions that facilitate scalability with low computational costs. As a step towards this goal, we pose the question: can learned policies be transferred across systems with higher number of agents without any extra training necessary and loss of performance?

Centralization and decentralization are two main paradigms in cooperative MARL. Centralization aims to reach cooperation naturally by transforming the multi-agent problem into a single agent problem by aggregating the local observations of agents to form a global state and inferring actions of agents through the global state. However, besides its inherent problems, full centralization with conventional function approximators does not allow execution on systems with different number of agents than the number of agents used in training, which we call transfer across variable number of agents. Nevertheless, use of graph convolutional networks (Kipf and Welling 2017), (Veličković et al. 2018) in a centralized training centralized execution paradigm as in (Li et al. 2020) allows transfer to systems with variable number of agents.

Centralized training decentralized execution paradigm makes transfer across variable number of agents straightforward, since the obtained policy infers actions from local observations and shares parameters with every agent in the system. Apart from that advantage, the overhead from increasing number of agents during execution is relatively small with respect to centralized execution. Therefore, we adopt centralized training decentralized execution as our main approach. Furthermore, we base our approach on graph convolutional MARL methods such as in (Jiang et al. 2020), (Li et al. 2020) that represent the locality and strength of agent interactions and form an implicit coordination graph.

In order to evaluate the generalization and transfer capacity across variable number of agents, we perform training in predator prey and traffic junction environments with varying number of agents starting from two-three agents to the number of agents that the environment's capacity allows. We then perform greedy evaluations on the trained models with varying number of agents. We analyze the performance of evaluation results in terms of generalization and transfer capacity for each combination of agent count for training versus evaluation. The experiment results demonstrate the possibility of obtaining similar or higher evaluation performance by training with less agents. We argue that training with a smaller number of agents and then transferring the model to high-scale configurations can be a more efficient solution than training with the high-scale configurations for agent count while preserving performance.

## Related Work

Centralization and decentralization are central approaches in cooperative MARL. Centralization aims to reach cooperation naturally by transforming the multi-agent problem into a single agent problem. However, (Sunehag et al. 2017) demonstrate that centralization leads to inefficient policies due to the lazy agent problem where some agents refrain from learning while an agent learns successfully because those agents' exploratory behaviour would damage the agent's learning performance. Furthermore, achieving centralization is not practically possible for some tasks due to impossibility or impracticability of being informed of other agents' observations. The alternative approach, and the mandatory approach when centralization is not possible, is decentralization where each agent is an independent learner. However, (Sunehag et al. 2017) assert that nonstationarity introduces spurious reward signals that an agent can not determine if the signal is the outcome of its own action or other agents' actions, thus leading to failure. In order to mitigate these drawbacks, it is a common approach to use the centralized training decentralized execution paradigm. Counterfactual Multi-Agent Policy Gradients (COMA) is a classic example of this paradigm which is an actor critic algorithm with a centralized actor and a decentralized critic (Foerster et al. 2017). We also use an actor-critic algorithm, Proximal Policy Optimization (PPO), as our base algorithm. PPO (Schulman et al. 2017) is a policy gradient method that optimizes a surrogate objective function that enables the algorithm to learn while limiting the extent the policy may change in each iteration.

Instead of adopting either centralized or decentralized approach, (Guestrin, Koller, and Parr 2002) propose formulating the cooperation problem by forming coordination graph and transforming this graph into a Dynamic Bayesian Network (DBN) for factorized representation, which allows factoring value functions in order to enable agents to coordinate by message passing and solved by linear programming. (Böhmer, Kurin, and Whiteson 2020) applies coordination graph approach to deep neural networks and approximate pay-off functions using them while maximizing value function by message passing.

Considering the convenience of representing multi-agent dynamics as a graph, processing graph structured data using deep neural networks is of importance. (Kipf and Welling 2017) propose a layer-wise propagation rule for deep neural networks processing graph structured data. (Veličković et al. 2018) introduce Graph Attention Networks which improve upon previous methods by using a masked self-attention layer that weights the impact of each neighbor accordingly during aggregation. (Jiang et al. 2020) propose using graph convolution with relation kernels in MARL to capture agent interplay that adapts to underlying dynamic graph of the environment in order to promote cooperation. (Li et al. 2020) suggest using self attention to obtain the coordination graph structure once and then using it for graph convolution in each pass to form an implicit deep coordination graph.

Our work adopts centralized training decentralized execution paradigm that uses PPO. We use graph convolution with self attention to process graph structured data and en-
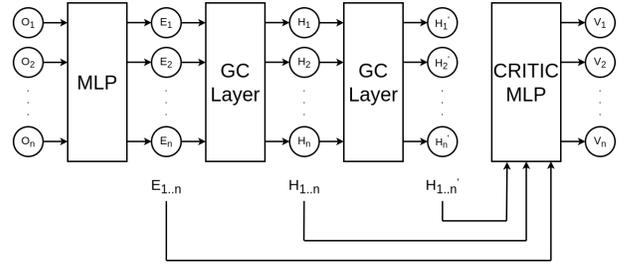


Figure 1: Critic Network - $o_i$ denotes observation of agent $i$, $E_i$ denotes embedding of the observation of agent $i$, $H_i$ denotes output of agent $i$ from the first graph convolution, $H_i'$ denotes output of agent $i$ from the second graph convolution, $v_i$ denotes the corresponding value of agent $i$

courage cooperation. Although the aforementioned related work also aims to promote cooperation, they mainly focus on achieving it for a fixed number of agents. Our contribution is checking the limits of cooperation that is learned for a fixed number of agents across different number of agents, and demonstrating that increasing agent count beyond a threshold in training is not necessary for achieving cooperation.

## Methodology

### Algorithm and Network Architecture

We adopt the Decentralized Partially Observable Markov Decision Process (Dec-POMDP) framework as formulated in (Oliehoek 2012) for this MARL problem. We use PPO with clipped surrogate objective as used and parametrized in (Li et al. 2020). We use a simple Multi-Layer Perceptron (MLP) that takes single local observations and outputs action probabilities for each agent.

The network architecture of centralized critic is given in Figure 1. The architecture consists of one MLP for extracting embeddings of decentralized observations. Then these embeddings are forwarded to two layers of graph convolutional layers with self attention which use message passing to aggregate these local observations properly. Then both embeddings and outputs of graph convolutional layers are forwarded to Critic MLP module which first aggregates information using sum operation and forwards it through MLP.

The network architecture of graph convolutional layer is provided in Figure 2. This layer passes each of its inputs through self attention, then it forwards attention outputs and residuals to a one layer neural network which produces graph convolutional layer outputs. The used self attention formulation and graph convolution outputs are inspired from (Jiang et al. 2020) and formulated in Equation 1 and Equation 2. In these equations, h denotes input to the graph convolutional layer, $h'$ denotes output of the graph convolutional layer, $d_k$ denotes the scaling factor, $\sigma$ denotes one layer feed forward network, and $W_Q, W_K, W_V$ denote weight matrices of query, key, and value vectors.
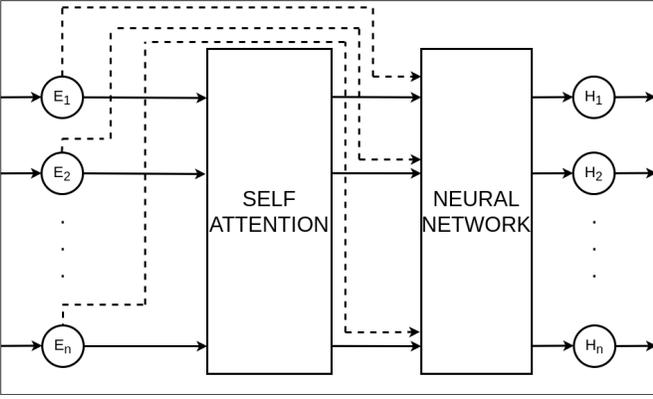
Figure 2: Graph Convolutional Layer - $E_i$ denotes the input of agent $i$ to the graph convolution, $H_i$ denotes the output of agent from the graph convolution

$$\text{Self-Attention}(h) = \text{softmax}\left(\frac{\mathbf{W}_Q h (\mathbf{W}_K h)^T}{\sqrt{d_k}}\right)\mathbf{W}_V h \tag{1}$$

$$h' = \sigma\left(\text{ concatenate }[\text{Self-Attention}(h), h]\right) \tag{2}$$

## Generalization and Transfer Capacity Evaluation

In order to evaluate the generalization and transfer capacity of the agents, we determine a range of agent count that starts from 2-3 and reaches to the capacity of the environment in question. Since we use decentralized execution that shares the actor with all agents, transfer is achieved simply by using this actor for all the agents with their local observations. We perform training in our chosen MARL environments for the whole range of agent count with 3 different seeds. We then determine a range of agent count for evaluation. The evaluation agent count range coincides with the training agent count range except for some additional sample points we may choose from inside the training range. The reason for choosing additional sample points from inside the training range is to make sure that we can answer the question: would training a model with higher agent count give superior results? Each evaluation process takes 100 evaluations averaged over with 3 different seeds in the environment, thus every combination contains 9 different evaluation results. At the end, evaluations of models with the same agent count for training and same agent count for evaluation are averaged and their standard deviations are calculated. For analysis, evaluation results are grouped according to the agent count during evaluation, in order to answer the question: is training $n$ agents the optimal choice for applications where $n$ or higher number of agents will be necessary, or is it possible to train with fewer number of agents and transfer the model to applications with higher number of agents necessary with similar or better performance?

```
input  : number of training epochs (ep), train agent
         count (tac), eval agent count (eac), train
         seeds, evaluation seeds
output: A
// A: 2d array with shape eac, tac
for each n in tac do
    for each seed in train seeds do
        for epoch ← 1 to to ep do
         │  Train agent_n for tac n
        end
        save agent_n with train seed: seed and tac: n
    end
end
for each n_eval in eac do
    for each n_train tac do
        result_eval ← 0
        for each agent in saved agents with train
          agent count = n_train do
            for each seed in evaluation seeds do
                result_avg ← evaluate 100 times
                result_eval ← result_eval + result_avg
            end
        end
        result_eval ← result_eval / (# of train seeds ×
          # of evaluation seeds)
        A[n_eval][n_train] ← result_eval
    end
end
```

**Algorithm 1:** Transfer Capacity Evaluation

# Experiments

## Predator Prey

Predator Prey environment consists of preys and predators where predators get rewarded by hunting preys, which move by hard coded action descriptions and random moves as described in (Li et al. 2020). As applied by (Li et al. 2020), we also penalize single agent capture attempts by -0.5 penalty in order to compel predators to collaborate. We use predator prey environment with a grid size of $20 \times 20$ in order to create capacity for 80 preys and 80 predators. We determine the training range and evaluation range of agent count as: 2, 5, 10, 20, 50, 80. The mean evaluation rewards are provided in Table 1 and Table 2.

Results of the experimental evaluation in predator prey

|    | 2               | 5               | 10              |
|----|-----------------|-----------------|-----------------|
| 2  | $-2.01 \pm 1.56$  | $30.30 \pm 1.95$  | $86.48 \pm 2.69$  |
| 5  | $+4.62 \pm 1.46$  | $41.98 \pm 0.77$  | $94.60 \pm 0.78$  |
| 10 | $-2.93 \pm 1.10$  | $39.71 \pm 1.56$  | $95.35 \pm 0.21$  |
| 20 | $-11.75 \pm 2.09$ | $22.05 \pm 6.95$  | $84.60 \pm 7.86$  |
| 50 | $-16.52 \pm 1.68$ | $6.68 \pm 4.03$   | $63.39 \pm 4.61$  |
| 80 | $-13.96 \pm 1.00$ | $11.17 \pm 4.29$  | $62.57 \pm 7.94$  |

Table 1: Mean of Total Rewards for Predator Prey (Columns denote the number of agents in evaluation while rows the denote number of agents in training.)

|     | 20 | 50 | 80 |
| --- | --- | --- | --- |
| 2 | $192.38 \pm 1.88$ | $496.58 \pm 0.56$ | $797.30 \pm 0.59$ |
| 5 | $196.37 \pm 0.58$ | $497.96 \pm 0.40$ | $798.07 \pm 0.82$ |
| 10 | $196.95 \pm 0.07$ | $498.18 \pm 0.37$ | $798.28 \pm 0.76$ |
| 20 | $194.08 \pm 2.48$ | $496.80 \pm 0.47$ | $794.74 \pm 2.30$ |
| 50 | $182.00 \pm 2.86$ | $494.68 \pm 1.00$ | $795.90 \pm 1.45$ |
| 80 | $168.27 \pm 8.47$ | $484.42 \pm 4.33$ | $789.59 \pm 2.09$ |

Table 2: Mean of Total Rewards for Predator Prey (Columns denote the number of agents in evaluation while rows the denote number of agents in training.)
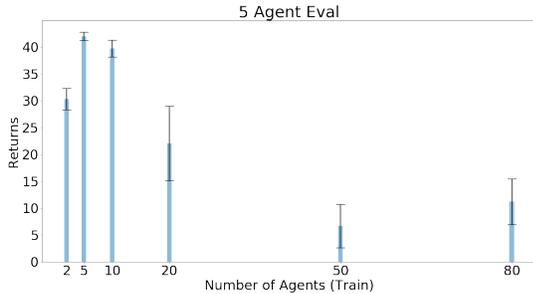
|     | 3 | 5 | 10 | 15 | 20 |
| --- | --- | --- | --- | --- | --- |
| 3 | $0.99 \pm 0$ | $0.92 \pm 0.04$ | $0.56 \pm 0.10$ | $0.26 \pm 0.14$ | $0.23 \pm 0.15$ |
| 5 | $0.99 \pm 0$ | $0.97 \pm 0.01$ | $0.77 \pm 0.09$ | $0.58 \pm 0.16$ | $0.58 \pm 0.18$ |
| 10 | $1.00 \pm 0$ | $0.99 \pm 0.00$ | $0.95 \pm 0.01$ | $0.84 \pm 0.07$ | $0.73 \pm 0.09$ |
| 15 | $1.00 \pm 0$ | $0.99 \pm 0.01$ | $0.94 \pm 0.01$ | $0.85 \pm 0.03$ | $0.79 \pm 0.05$ |
| 20 | $0.99 \pm 0$ | $0.99 \pm 0.00$ | $0.90 \pm 0.02$ | $0.83 \pm 0.04$ | $0.79 \pm 0.03$ |

Table 3: Mean of Success Rates for Traffic Junction (Columns denote the number of agents in evaluation while rows denote the number of agents in training.)
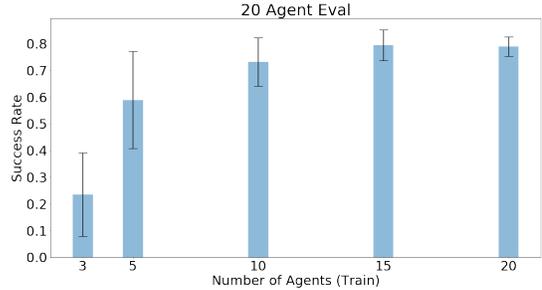


Figure 3: 5 Agents Evaluation in Predator Prey



Figure 4: 20 Agents Evaluation in Traffic Junction

environment, as provided in Table 1, Table 2, and Figure 3, demonstrate that training in few number of agents such as 2-5 can get evaluation results that surpass the evaluation results of models that are trained with large number of agents such as 50-80. It can be deduced that models trained with few number of agents have high generalization and transfer capacity for execution with high number of agents. However, our analysis shows us that the reverse is not true. The models that are trained with large number of agents such as 50-80 have very low returns for the evaluation cases where there are 2-5 agents in the environment. Hence, it can be inferred from the results that for a high performance application of predator prey environment, choosing number of agents to train from the range [5, 10] would be the better choice.

## Traffic Junction

Traffic Junction environment consists of predetermined routes and junctions where cars need to reach their destination point without collision as described in (Li et al. 2020). We use traffic junction environment with the difficulty 'hard'. Other parameters of the environment such as dimension, maximum agent add rate and minimum add agent rate are chosen compatibly with the environment difficulty according to the configurations proposed by (Singh, Jain, and Sukhbaatar 2018). Because of environment dimension and agent add rate setting, the capacity of the environment allows for approximately 20 agents in a single time step. Thus, we determine the training range and evaluation range of agent count as: 3, 5, 10, 15, 20. The mean evaluation success rates are provided in Table 3.

Results of the experimental evaluation in traffic junction environment, as provided in Table 3 and Figure 4, demonstrate that training with few number of agents such as 3-5

gets evaluation results with much lower success rate compared to the evaluation results of models that are trained with large number of agents such as 15-20. It can be deduced that models trained with few number of agents can not sufficiently transfer for execution with high number of agents. The models that are trained with 15-20 agents have very high success rate for the evaluation cases where there are 3-5 agents in the environment. Nevertheless, the evaluation results with 15 and 20 agents demonstrate that training with 15 agents gives better results than training with 20 agents. Thus, we infer that an environment can have a sweet spot for the number of agents to train. We also infer that environment dynamics play a key role in the generalization and transfer capacity of training.

## Discussion

In this work, we adopted the centralized training decentralized execution paradigm and investigated the generalization and transfer capacity of the trained models across variable number of agents. We assessed it by training variable number of agents in a specific MARL problem and then performing greedy evaluations with variable number of agents for each training configuration. We deduced that an environment can have a sweet spot for the number of agents to train in terms of evaluation performance and that environment dynamics play a key role in the generalization and transfer capacity of training. We saw that training with fewer number of agents can be a more efficient option for execution in large number of agents. We conclude that optimal number of agents to perform training may differ from the target number of agents and put forward that transfer across large number of agents can be a more efficient solution to scaling up than directly increasing number of agents during training.

## Acknowledgments

## References

Böhmer, W.; Kurin, V.; and Whiteson, S. 2020. Deep Coordination Graphs.

Foerster, J.; Farquhar, G.; Afouras, T.; Nardelli, N.; and Whiteson, S. 2017. Counterfactual Multi-Agent Policy Gradients.

Guestrin, C.; Koller, D.; and Parr, R. 2002. Multiagent Planning with Factored MDPs. In Dietterich, T. G.; Becker, S.; and Ghahramani, Z., eds., *Advances in Neural Information Processing Systems 14*, 1523–1530. MIT Press. URL http://papers.nips.cc/paper/1941-multiagent-planning-with-factored-mdps.pdf.

Jiang, J.; Dun, C.; Huang, T.; and Lu, Z. 2020. Graph Convolutional Reinforcement Learning.

Kipf, T. N.; and Welling, M. 2017. Semi-Supervised Classification with Graph Convolutional Networks.

Li, S.; Gupta, J. K.; Morales, P.; Allen, R.; and Kochenderfer, M. J. 2020. Deep Implicit Coordination Graphs for Multi-agent Reinforcement Learning.

Oliehoek, F. A. 2012. *Decentralized POMDPs*, 471–503. Berlin, Heidelberg: Springer Berlin Heidelberg. ISBN 978-3-642-27645-3. doi:10.1007/978-3-642-27645-3_15. URL https://doi.org/10.1007/978-3-642-27645-3_15.

Schulman, J.; Wolski, F.; Dhariwal, P.; Radford, A.; and Klimov, O. 2017. Proximal Policy Optimization Algorithms.

Singh, A.; Jain, T.; and Sukhbaatar, S. 2018. Learning when to Communicate at Scale in Multiagent Cooperative and Competitive Tasks.

Sunehag, P.; Lever, G.; Gruslys, A.; Czarnecki, W. M.; Zambaldi, V. F.; Jaderberg, M.; Lanctot, M.; Sonnerat, N.; Leibo, J. Z.; Tuyls, K.; and Graepel, T. 2017. Value-Decomposition Networks For Cooperative Multi-Agent Learning. *CoRR* abs/1706.05296. URL http://arxiv.org/abs/1706.05296.

Veličković, P.; Cucurull, G.; Casanova, A.; Romero, A.; Liò, P.; and Bengio, Y. 2018. Graph Attention Networks.