

Bayesian and Attentive Aggregation for Cooperative Multi-Agent Deep Reinforcement Learning

Master's Thesis
of

Robin Ruede

KIT Department of Informatics
Institute for Anthropomatics and Robotics (IAR)
Autonomous Learning Robots (ALR)

Referees: Prof. Dr. Techn. Gerhard Neumann
Prof. Dr. Ing. Tamim Asfour

Advisor: M.Sc. Maximilian Hüttenrauch

Duration: February 1st, 2021 — July 31st, 2021

Plagiarism Declaration (German)

Ich versichere hiermit, dass ich die Arbeit selbstständig verfasst habe, keine anderen als die angegebenen Quellen und Hilfsmittel benutzt habe, die wörtlich oder inhaltlich übernommenen Stellen als solche kenntlich gemacht habe und die Satzung des Karlsruher Instituts für Technologie zur Sicherung guter wissenschaftlicher Praxis beachtet habe.

Karlsruhe, den July 31st, 2021

Robin Ruede

German Summary

Mehr-Agenten-Reinforcement-Learning (MARL) ist ein aufkommendes Feld des Reinforcement Learning mit realen Anwendungen wie unbemannte Luftfahrzeuge, Such- und Rettungsdienste und Lagerorganisation. Es gibt viele verschiedene Ansätze, um die Methoden des Single-Agent Reinforcement Learning auf MARL anzuwenden. In dieser Arbeit geben wir einen Überblick über verschiedene Lernmethoden und Umgebungseigenschaften und konzentrieren uns dann auf ein Problem, das bei den meisten Varianten von MARL auftritt: Wie sollte ein Agent die Informationen nutzen, die er aus einer großen und variierenden Anzahl von Beobachtungen der Welt gesammelt hat, um Entscheidungen zu treffen? Wir konzentrieren uns auf drei verschiedene Methoden zur Aggregation von Beobachtungen und vergleichen sie hinsichtlich ihrer Trainingsleistung und Sample Efficiency. Wir stellen eine auf Bayes'scher Konditionierung basierende Policy-Architektur für die Aggregation vor und vergleichen sie mit der Mittelwertaggregation und der Attention-Aggregation, die in verwandten Arbeiten verwendet werden. Wir zeigen die Leistung der verschiedenen Methoden an einer Reihe von kooperativen Aufgaben, die auf eine große Anzahl von Agenten skaliert werden können, einschließlich Aufgaben, die andere Objekte in der Welt haben, die von den Agenten beobachtet werden müssen, um die Aufgabe zu lösen.

Wir optimieren die Hyperparameter, um zeigen zu können, welche Parameter zu den besten Ergebnissen für jede der Methoden führen. Darüber hinaus vergleichen wir verschiedene Varianten der Bayes'schen Aggregation und vergleichen die kürzlich eingeführte Lernmethode Trust Region Layers mit der allgemein bekannten Proximal Policy Optimization (PPO).

Abstract

Multi-agent reinforcement learning (MARL) is an emerging field in reinforcement learning with real world applications such as unmanned aerial vehicles, search-and-rescue, and warehouse organization. There are many different approaches for applying the methods used for single-agent reinforcement learning to MARL. In this work, we survey different learning methods and environment properties and then focus on a problem that persists through most variants of MARL: How should one agent use the information gathered from a large and varying number of observations of the world in order to make decisions? We focus on three different methods for aggregating observations and compare them regarding their training performance and sample efficiency. We introduce a policy architecture for aggregation based on Bayesian conditioning and compare it to mean aggregation and attentive aggregation used in related work. We show the performance of the different methods on a set of cooperative tasks that can scale to a large number of agents, including tasks that have other objects in the world that need to be observed by the agents in order to solve the task.

We optimize the hyperparameters to be able to show which parameters lead to the best results for each of the methods. In addition, we compare different variants of Bayesian aggregation and compare the recently introduced Trust Region Layers learning method to the commonly used Proximal Policy Optimization.

Table of Contents

German Summary	v
Abstract	vii
1. Introduction	1
2. Preliminaries	3
2.1. Reinforcement Learning	3
2.2. Actor-Critic Training Methods	5
2.2.1. Proximal Policy Optimization (PPO)	5
2.2.2. Trust Region Layers (PG-TRL)	6
2.3. Multi-Agent Reinforcement Learning (MARL)	6
2.4. Environment Model and Learning Process	7
2.5. Aggregation Methods	8
2.5.1. Concatenation	8
2.5.2. Mean Aggregation	9
2.5.3. Aggregation With Other Pooling Functions	9
2.5.4. Bayesian Aggregation	10
2.5.5. Attention Mechanisms	11
2.5.6. Other Aggregation Methods	12
3. Related Work	13
3.1. Centralized vs Decentralized Learning	13
3.2. Centralized vs Decentralized Execution	14
3.3. Cooperative, Adversarial, Team-Based	14
3.4. Partial Visibility	15
3.5. Simultaneous vs Turn-Based	15
3.6. MARL Tasks in Related Work	16

4. Scalable Information Aggregation for Deep MARL Policies	17
4.1. Policy Architecture	17
4.2. Mean/Max/Softmax Aggregation	19
4.3. Bayesian Aggregation	21
4.4. Attentive Aggregation	22
5. Experiments	25
5.1. Experimental Setup	25
5.2. Considered Tasks	26
5.2.1. Rendezvous Task	26
5.2.2. Single-Evader Pursuit Task	27
5.2.3. Multi-Evader Pursuit	29
5.2.4. Box Assembly Task	29
5.2.5. Box Clustering Task	30
6. Results	32
6.1. Aggregation Method Results	33
6.1.1. Rendezvous Task	33
6.1.2. Single-Evader Pursuit Task	34
6.1.3. Multi-Evader Pursuit Task	36
6.1.4. Assembly Task	40
6.1.5. Clustering Task With Two Clusters	40
6.1.6. Clustering Task With Three Clusters	40
6.2. Learning Algorithm Comparison (PPO vs PG-TRL)	42
6.3. Bayesian Aggregation Variants	44
6.3.1. Single Space vs Separate Space Aggregation	44
6.3.2. Separate vs Common Encoder	44
6.3.3. Using the Aggregated Variance or Only the Mean	46
7. Conclusion and Future Work	48
7.1. Conclusion	48
7.2. Future Work	49
Bibliography	51
A. Experiment Hyper-Parameters and Overview	57

Chapter 1.

Introduction

Ant colonies, bee swarms, fish colonies, and migrating birds all exhibit swarming behavior to achieve a common goal or to gain an advantage over what's possible alone. Each individual within a swarm usually only has limited perception, and often there is no central control that enforces a common goal or gives instructions. Swarms of animals can self-organize and exhibit complex emergent behavior in spite of the limited intelligence, limited perception, and limited strength of every individual.

In recent years and with the advent of deep reinforcement learning, it has become possible to create artificial agents that solve fairly complex tasks in simulated environments as well as in the real world, even when the path to the goal is difficult to identify and the reward is sparse. Most of the research, however, is focused on a single agent interacting with a world, and giving the single agent as much power and flexibility as it needs to solve the task. Akin to animal swarms, which often consist of fairly simple and limited individuals, we focus on simple artificial agents that need to cooperate to solve a common task, where the task is either difficult or impossible for an individual agent. The agents thus need to be able to learn to work together and to act even with the limited information they are able to perceive.

Swarms of artificial agents can have multiple applications, some inspired by swarms in the animal kingdom, some going beyond. Robots can be used in logistics to organize warehouses. Swarms of robots can be used to map out dangerous areas, solve mazes, find trapped or injured people in search-and-rescue missions quicker and more safely than humans [Dre21]. Swarms of walking robots could be used for animal shepherding [Hu+20]. Drone swarms can be used for entertainment in light shows [WKA17], for advertising, for autonomous surveillance [Bür09], or they could be weaponized and

used as military robots [San17]. Further in the future, swarms of nanobots could be used in medicine to find and target specific cells in the body [CS16]. Swarms of autonomous spacecraft can be used to explore space, to harvest resources, and to terraform planets into paperclips [TD20].

Swarms of homogenous agents can be very resilient. Since they can be designed such that no one individual is a critical component of the swarm as a whole, individual failures do not necessarily result in a critical collapse of the whole swarm.

Robots can be controlled by explicit algorithmic behaviour, but that can be complicated, inflexible and fragile, as is shown by the success in recent applications of reinforcement learning to control tasks that have not been solved by pre-programmed algorithms [Ope+19]. Controlling robot swarms using policies learned with deep reinforcement learning shows promising results in recent literature such as [Bak+20].

To apply deep reinforcement learning to multi-agent systems, we need to make a few adjustments to the existing learning algorithms and figure out how best to design the policy network. Specifically, we need a way to feed a large and varying amount of observations from the neighboring agents into the fixed size input of the dense neural network representing the policy. In this work, we consider three aggregation methods on a set of different multi-agent tasks: Mean aggregation, Bayesian aggregation, and attentive aggregation. Our main goal is to compare these methods with regard to their training performance and sample efficiency. We limit ourselves to a specific subset of MARL tasks that are fully cooperative with a team reward and with homogenous agents in a centralized-learning/decentralized-execution setup.

We first give an overview over all the preliminaries we need for our work in Chapter 2, including reinforcement learning in general, the multi-agent extensions for reinforcement learning, and the background for the aggregation methods we use. Next, we describe some related work in Chapter 3. Then, we describe our contribution in Chapter 4 with details about the policy architecture and the specifics for the different aggregation methods. Our experimental setup, including the specific environments we use to carry out our experiments are described in Chapter 5. Finally, we show and interpret the results of our experiments in Chapter 6 and talk about the conclusions we can draw from the experiments as well as the potential for future work in Chapter 7.

Chapter 2.

Preliminaries

In this chapter, we introduce the preliminaries we need for our work. We describe reinforcement learning and two specific policy gradient training methods used for deep reinforcement learning in general and then the specifics of RL for multi-agent systems, including the need for information aggregation and the different methods we consider.

2.1. Reinforcement Learning

Reinforcement learning is a method of training an artificial agent to solve a task in an environment. As opposed to supervised learning, there is no explicit label / path to a solution. Instead, the agent iteratively takes actions that affect its environment, and receives a reward when specific conditions are met. The reward can be dense (positive or negative signal at every step) or very sparse (only a binary reward at the end of an episode).

Reinforcement learning problems are usually defined as Markov Decision Processes (MDPs). An MDP is an extension of a Markov chain, which is a set of states with a transition probability between each pair of states. The probability only depends on the current state, not the previous states.

MDPs extend Markov chains, adding a set of actions (that allow an agent to influence the transition to the next state) and rewards (that motivate the agent). An MDP thus consists of:

- A set of states S (the *state space*)
- A set of actions A (the *action space*)
- The transition probability that a specific action in a specific state leads to specific second state $T(s, a, s') = P(s'|s, a)$
- The reward function that specifies the immediate reward an agent receives depending on the previous state, the action, and the new state $R: S \times A \times S \rightarrow \mathbb{R}$

The MDP can either run for an infinite period or finish after a number of transitions. The transition function then always returns the same state.

The method by which an agent chooses its actions based on the current state is called a *policy* π . The policy defines a probability for taking each action based on a specific state.

Based on a policy we can also define the *Q-value function* and the *value function*. The value function V^π is the sum of all future expected rewards that an agent gets based on an initial state s_t and a specific policy that leads to expected rewards r in each future time step $t+i$ when sampling actions according to the policy:

$$V^\pi(s_t) = \mathbb{E} \left[\sum_{i=0}^{\infty} \gamma^i r_{t+i} \right]$$

The value function can include an exponential decay γ for weighing future rewards, which is especially useful if the horizon is infinite. The Q-value function is the same except that the action at step t is fixed to a_t instead of chosen by the policy:

$$Q^\pi(s_t, a_t) = \mathbb{E} \left[\sum_{i=0}^{\infty} \gamma^i r_{t+i} \right]$$

Often we need to extend MDPs to allow for an observation model: A *partially observable Markov decision process* (POMDP) is an extension to MDPs that introduces an indirection between the set of states and the input to the policy (called an *observation*). In the definition of a POMDP a set of observations is added with a probability of each state leading to a specific observation. The policy now depends on the observation and the action instead of the state and the action. An observation does not necessarily include all the information from the corresponding state, which leads to naming “partially observable”.

To successfully solve a reinforcement learning task, we need to find a policy that has a high expected reward — we want to find the *optimal policy function* that has the highest *value* on the initial states of our environment. Since finding the optimal policy

directly is impossible for even slightly complicated tasks, we instead use optimization techniques.

2.2. Actor-Critic Training Methods

Policy gradient training methods are a reinforcement learning technique that optimize the parameters of a policy using gradient descent. Actor-critic methods optimize both the policy and an estimation of the value function at the same time.

There are multiple commonly used actor critic training methods such as Trust Region Policy Optimization (TRPO) [Sch+15], Proximal Policy Optimization (PPO) [Sch+17] and Soft Actor Critic [Haa+18]. We run most of our experiments with one of the most commonly used methods (PPO) and also explore a new method that promises stabler training (PG-TRL [Ott+20]).

2.2.1. Proximal Policy Optimization (PPO)

Proximal Policy Optimization [Sch+17] is an actor-critic policy gradient method for reinforcement learning. In PPO, each training step consists of collecting a set of trajectory rollouts, then optimizing a “surrogate” objective function using stochastic gradient descent. The surrogate objective function approximates the policy gradient while enforcing a trust region by clipping the update steps. PPO is a successor to Trust Region Policy Optimization (TRPO) with a simpler implementation and empirically better performance.

PPO optimizes the policy using

$$\theta_{k+1} = \operatorname{argmax}_{\theta} E_{s,a \sim \pi_{\theta_k}} [L(s, a, \theta_k, \theta)]$$

Where π_{θ_k} is a policy with parameters θ in training step k , s is the state, $a \sim \pi_{\theta_k}$ is the action distribution according to the policy at step k , and $\varepsilon = 0.2$ is a hyperparameter. L is given by

$$L(s, a, \theta_k, \theta) = \min \left(\frac{\pi_{\theta}(a|s)}{\pi_{\theta_k}(a|s)} A^{\pi_{\theta_k}}(s, a), \operatorname{clip}\left(\frac{\pi_{\theta}(a|s)}{\pi_{\theta_k}(a|s)}, 1 - \varepsilon, 1 + \varepsilon\right) A^{\pi_{\theta_k}}(s, a) \right).$$

The advantage A is defined as $A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$ and is estimated using Generalized Advantage Estimation [Sch+18] based on the estimated value function.

We use PPO as the default for most of our experiments since it is widely used in other deep reinforcement learning work.

2.2.2. Trust Region Layers (PG-TRL)

Differentiable trust region layers are an alternative method to enforce a trust region during policy updates introduced by Otto et al. [Ott+20]. PPO uses a fixed clipping ratio to enforce the trust region, which can result in unstable training. Instead of using a fixed clipping boundary, in PG-TRL a surrogate policy is created with a new layer at the end that projects the unclipped policy back into the trust region. The trust region and the projection is based on either the KL-divergence [KL51] or the Wasserstein W_2 distance [Vil09].

After each training step, the projected new policy depends on the previous policy. To prevent an infinite stacking of old policies, the explicitly projected policy is only used as a surrogate, while the real policy is based on a learned approximation. To prevent the real policy and the projected policy from diverging, the divergence between them is computed again in the form of the KL divergence or the Wasserstein W_2 distance. The computed divergence (*trust region regression loss*) is then added to the policy gradient loss function with a high factor.

We explore PG-TRL as an alternative training method to PPO.

2.3. Multi-Agent Reinforcement Learning (MARL)

Usually, reinforcement learning environments are modeled as POMDPs. POMDPs only have a single agent interacting with the environment, so for multi-agent settings, we need a different system that can model multiple agents interacting with the world. There are different ways of extending POMDPs to work for multi-agent tasks. Decentralized POMDPs (Dec-POMDP) [OA16] are a generalization of POMDPs that split the action and observation spaces into those of each agent.

A Dec-POMDP consists of:

- A set of n agents $D = \{1, \dots, n\}$
- A set of states S
- A set of joint actions A
- A set of transition probabilities between states $T(s, a, s') = P(s' | s, a)$

- A set of joint observations O
- An immediate reward function $R : S \times A \rightarrow \mathbb{R}$

Note that this definition is very similar to a normal POMDP. The difference is that the set of joint observations and joint actions consists of a tuple of the observations/actions of each individual agent (i.e. $a \in A = \langle a_1, \dots, a_n \rangle$, where a_i is the action of agent i).

In Dec-POMDPs every agent can have differing observations and actions. The subset of Dec-POMDPs where the agents are homogenous has been described as a SwarMDP by Šošić et al. [Šoš+17]. The joint observation and action spaces thus become $O = \hat{O}^n$ and $A = \hat{A}^n$, where \hat{O} and \hat{A} are the observation and action space of each agent.

2.4. Environment Model and Learning Process

In our experiments, we impose a set of restrictions on the environments and learning process. The restrictions we impose are mostly based on [HŠN19]. A more detailed description of these restrictions, as well as related work using other variants is given in in Chapter 3.

In general, the agents in a multi-agent environment can differ in their intrinsic properties. For example, they can have different control dynamics, maximum speeds, different observation systems, or different possible actions. We only consider environments with homogenous agents: All agents have the same physical properties, observation space, and action space. They only differ in their extrinsic properties, e.g., their current position, rotation, and speed. This also causes them to have a different perspective, different observations and thus different actions, resulting in differing behavior even when they are acting according to the same policy. We thus only consider environments where the agents are homogenous — each agent has the same possible observations and actions.

We only consider cooperative environments, and we use the same reward function for all agents. Real-world multi-agent tasks are usually cooperative since in adversarial environments, where different agents have different goals, one entity would not have control over multiple adversarial parties.

We focus on global visibility since the additional noise introduced by local observability would be detrimental to the quality of our results.

For training, we use the centralized-learning/decentralized-execution (CLDE) approach — a shared common policy is learned for all agents, but the policy is executed by each agent separately.

PPO and other policy gradient methods are designed for single-agent environments. We adapt PPO to the multi-agent setting without making any major changes: The policy parameters are shared for all agents, and each agent gets the same global reward. The value function for advantage estimation is based on the same architecture as the policy. Since the stable-baselines3 implementation of PPO is already written for vectorized environments (collecting trajectories from many environments running in parallel), we create a new VecEnv implementation that flattens multiple agents in multiple environments.

Similarly to the setup used for TRPO by Hüttenrauch, Šošić, and Neumann [HŠN19], we collect the data of each agent as if that agent was the only agent in the world. For example, a batch of a single step of 10 agents each in 20 different environment becomes 200 separate training samples. Each agent still only has access to its own local observations, not the global system state. This means that during inference time, each agent has to act independently, based on the observations it makes locally, but is rewarded with the same reward each agent receives. This common reward reflects the cooperative structure of the task.

2.5. Aggregation Methods

The observations of each agent in an MARL task contains a varying number of observables. The observables can be clustered into groups where each observable is of the same kind and shape. For example, one observable group would contain all the neighboring agents, while another would contain, e.g., a specific type of other observed objects in the world.

To apply actor-critic training methods, both our policy to predict the actions and the critic to estimate the value function are trainable neural networks. We need a method to aggregate the observations in one observable group into a format that can be used as the input of these neural networks. Specifically, this means that the output of the aggregation method needs to have a fixed dimensionality. In the following, we present different aggregation methods and their properties, including related work that uses those aggregation methods.

2.5.1. Concatenation

The simplest aggregation method is concatenation, where each observable is concatenated along the feature dimension into a single observation vector. This method has a few issues however: We can have a varying number of observables, but since the input size of the neural policy network is fixed, we need to set a maximum number of observables, and set the input dimensionality of the policy architecture proportional

to that maximum. Concatenation also ignores the other useful properties of the observables, namely the uniformity (the feature at index i of one observable has the same meaning as the feature at index i of every other observable in a group) and the exchangeability (the order of the observables is either meaningless or variable). For concatenation, we have to choose an ordering of the observables. This ordering is either stable but meaningless (since the observables are by definition permutation-invariant), or meaningful but unstable (ordering by an extrinsic property, for example ordering by distance to the current agent).

Concatenation scales poorly with a large number of observables since the input size of the neural network has to scale proportionally to the maximum number of observables. It is used for example by Lowe et al. [Low+17] to aggregate the neighboring agents' observations and actions.

2.5.2. Mean Aggregation

Instead of concatenating each element o_i in an observable group O , we can also interpret each element as a sample of a distribution that describes the current system state. We use the empirical mean of the samples to retrieve a representation of the system state ψ_O based on all observed observables, as shown in Equation (2.1).

$$\psi_O = \mu_O = \frac{1}{|O|} \sum_{o_i \in O} o_i \quad (2.1)$$

Hüttenrauch, Šošić, and Neumann [HŠN19] used mean aggregation for deep multi-agent reinforcement learning and compared it to other aggregation methods. Gebhardt, Hüttenrauch, and Neumann [GHN19] applied mean aggregation to more complex tasks in more realistic simulated environments.

Mean aggregation is strongly related to mean field theory. Mean field theory is a general principle of modeling the effect that many particles have by averaging them into a single field, ignoring the individual variances of each particle. The application of mean field theory for multi-agent systems were formally defined by Lasry and Lions [LL07] as *Mean Field Games*. In MARL, mean field Q-learning and mean field actor-critic was defined and evaluated by Yang et al. [Yan+18]. Chen et al. [Che+21] use mean fields productively for control of numerous unmanned aerial vehicles.

2.5.3. Aggregation With Other Pooling Functions

Instead of using the empirical mean, other aggregation methods can also be used:

Max pooling:

$$\psi_O = \max_{o_i \in O} o_i$$

Softmax pooling:

$$\psi_O = \sum_{o_i \in O} o_i \frac{e^{o_i}}{\sum_{o_j \in O} e^{o_j}}$$

Max-pooling is widely used in convolutional neural networks to reduce the feature map dimensionality, and it was used as an alternative to mean pooling in MARL by [GHN19]. Softmax aggregation was used by Mordatch and Abbeel [MA18] for MARL.

2.5.4. Bayesian Aggregation

Aggregation with Gaussian conditioning works by starting from a Gaussian prior distribution and updating it using a probabilistic observation model for every seen observation. Gaussian conditioning is widely used in applications such as Gaussian process regression [Ras04].

We consider Bayesian aggregation as defined by [Vol+20, sec 7.1]. We introduce z as a random variable with a Gaussian distribution:

$$z \sim \mathcal{N}(\mu_z, \sigma_z^2) \quad (p(z) \equiv \mathcal{N}(\mu_z, \sigma_z^2))$$

Initially, this random variable is estimated using a diagonal Gaussian prior as an a-priori estimate:

$$p_0(z) \equiv \mathcal{N}(\mu_{z_0}, \text{diag}(\sigma_{z_0}^2))$$

This prior is then updated with Bayesian conditioning using each of the observed elements r_i . We interpret each observation as a new sample from the distribution $p(z)$, each with a mean r_i and a standard deviation σ_{r_i} . We use the probabilistic observation model and consider the conditional probability

$$p(r_i|z) \equiv \mathcal{N}(r_i|z, \sigma_{r_i}^2).$$

With standard Gaussian conditioning [Bis06] this leads to the closed form factorized posterior description of z :

$$\sigma_z^2 = \frac{1}{\frac{1}{\sigma_{z_0}^2} + \sum_{i=1}^n \frac{1}{\sigma_{r_i}^2}}$$

$$\mu_z = \mu_{z_0} + \sigma_z^2 \cdot \sum_{i=1}^n \frac{(r_i - \mu_{z_0})}{\sigma_{r_i}^2}$$

Bayesian aggregation was used by Volpp et al. [Vol+20] for context aggregation in conditional latent variable (CLV) models. There is no related work using Bayesian aggregation for MARL.

2.5.5. Attention Mechanisms

Attention mechanisms are commonly used in other areas of machine learning.

An attention function computes some compatibility between a *query* Q and a *key* K , and uses this compatibility as the weight to compute a weighted sum of the *values* V . The query, key, and value are all vectors.

The multi-head attention mechanism was introduced by Vaswani et al. [Vas+17] and is the core of the state-of-the-art model for many natural language processing tasks.

We only consider the specific multi-head residual masked self attention variant of attention mechanisms for observation aggregation used by Baker et al. [Bak+20].

There, the attention function in is a scaled dot-product attention as described by [Vas+17]:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

d_k is the dimensionality of the key K .

Instead of using only a single attention function, [Vas+17] uses multiple independent attention heads. The inputs (Q, K, V) of each of the heads $1, \dots, n$ as well as the concatenated output is transformed with separately learned dense layers (described as weight matrices W_i^Q, W_i^K, W_i^V, W^O). The full multi-head attention module $\text{MHA}()$ thus looks like this:

$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

$$\text{MHA}(Q, K, V) = \text{concat}(\text{head}_1, \dots, \text{head}_n)W^O$$

For self-attention, the query Q , the key K and the value V are all set to the same input value. In [Bak+20], the authors combine the multi-head attention with a mean aggregation. Note that they only use the attention mechanism to individually transform the information from each separate observable into a new feature set instead of directly using it as a weighing function for the aggregation.

Iqbal and Sha [IS19] use a different approach with a decentralized policy and a centralized critic. An attention mechanism is used to weigh the features from the other agents, then aggregates them to calculate the value function.

Liu and Tan [LT20] first encode all aggregatable observations together with the proprioceptive observations with an “extrinsic encoder”, then compute attention weights using another encoder from that and aggregate the features retrieved from a separate “intrinsic encoder” using those weights.

2.5.6. Other Aggregation Methods

Hüttenrauch, Šošić, and Neumann [HŠN19] also did experiments with aggregation into histograms and aggregation with radial basis functions, though the results indicated they were outperformed by a neural network encoder with mean aggregation.

Chapter 3.

Related Work

There are many variants of applying reinforcement learning to multi-agent systems. In this chapter, we describe some related work in MARL, including variations in the training process and the environments.

An overview over recent MARL work and some of the differing properties can be found in [ZYB21] and [Can+21].

3.1. Centralized vs Decentralized Learning

During training, the agent policies can either be learned in a centralized or a decentralized fashion.

In decentralized learning, each agent learns their own policy, while in centralized learning, the policy is shared between all agents. Decentralized learning has the advantage that the agents do not need to be homogenous, since the learned policy can be completely different. It also means that it's harder for the agents to learn to collaborate though, since they are not intrinsically similar to each other. Decentralized learning is used in [Wai+19; KMP13; Zha+18].

Centralized learning means the training process happens centralized, with a single policy. It has the advantage of only needing to train one policy network and the possibility of being more sample-efficient. Centralized learning requires homogenous

agents since the policy network parameters are shared across all agents. It is used for example in [Foe+17; GEK17].

3.2. Centralized vs Decentralized Execution

When using centralized learning, it is possible to use a single policy to output actions for all agents at the same time. This is called “centralized execution”. The policy chooses the actions based on the global system state in a “birds-eye” view of the world. The disadvantage is that this means agents can’t act independently if the environment has restrictions such as partial observability, and it is less robust to communication failures and to agent “deaths”.

The alternative is decentralized execution — the observation inputs and action outputs of the policy networks are local to a single agent. This can formally be described as a *Decentralized Partially Observable Markov Decision Process* (Dec-POMDP) [Oli12]. Decentralized learning with decentralized execution is needed for heterogeneous agents (e.g. [Zha+18; KMP13]). When combining centralized learning with decentralized execution, there is only one policy network that is shared between all agents but the inputs and outputs of the policy are separate.

Centralized-learning with decentralized-execution is thus a compromise between performance, robustness, and sample efficiency. CLDE is used e.g. by [Low+17; HŠN19; GEK17; Foe+17; Foe+16].

3.3. Cooperative, Adversarial, Team-Based

Multi-agent environments can be cooperative, adversarial, or team-based. Cooperative environments are those where all the agents have one common goal. In adversarial tasks each agent has their own independent goal that conflicts with the goal of the other agents. An example for a cooperative environment is the rendezvous task: All agents need to meet up at a single point, where the agents have to decide independently on the location. The reward here is the negative average pairwise distance of the agents, see Section 5.2.1. Note that cooperative environments can also include other adversarial agents, as long as those agents are controlled by an explicit algorithm and not a learned policy. From the perspective of the policy, these agents are considered part of the environment.

Cooperative learning can be done with separate agent rewards or a single common reward. Cooperative learning with a single reward means the reward must be somewhat sparse, since each agent cannot easily judge what the impact of its own actions were on the reward in any given time step. This also makes it impossible for an agent

to gain an egoistic advantage over the other agents, enforcing the learned policy to become Pareto optimal. Another approach was introduced by Matignon, Laurent, and Le Fort-Piat [MLL07], who create a compromise between the sample-efficiency of individual rewards and the Pareto-optimality of a common reward for cooperative MARL settings with their *Hysteretic Q-Learning* algorithm that jointly optimizes an individual as well as a common reward.

Adversarial environments are usually zero-sum, that is the average reward over all agents is zero. An example for an adversarial environment is the *Gather* task described by Zheng et al. [Zhe+18]: In a world with limited food, agents need to gather food to survive. They can also kill each other to reduce the scarcity of the food. Another example of an adversarial task is Go [Sil+17] as well as many other board games, though these usually have a fairly small number of agents.

Team-based tasks have multiple teams that each have a conflicting goal, but from the perspective of each team the task is cooperative (with a single reward function). The team reward can either be defined directly for the team or by averaging a reward function of each team member.

An example of a team-based environment is the OpenAI Hide-and-Seek Task [Bak+20]. In this task, there are two teams of agents in a simulated physical world with obstacles, and members of one team try to find the members of the other team. The Hide-team is rewarded +1 if none of the team members is seen by any seeker, and -1 otherwise. The Seek-team is given the opposite reward.

3.4. Partial Visibility

The observations that each agent receives in our experiments are local. For example, if one agent sees another, that agent's properties are observed relative to the current agent — the distance, relative bearing, and relative speed.

In addition, each agent may only have local visibility, for example it can only observe the positions of agents and objects in the world within some radius or the visibility can be hindered by obstacles. In [HŠN19] both the local and global visibility variants of the same tasks were considered.

3.5. Simultaneous vs Turn-Based

In general, multi-agent environments can be turn-based or simultaneous. In turn-based environments each agent acts in sequence, with the world state changing after

each turn. In simultaneous environments, all agents act “at the same time”, i.e. the environment only changes after all agents have acted in one time step.

Simultaneous environments can be considered a subset of turn-based environments, since a simultaneous environment can be converted to a turn-based one by fixing the environment during turns and only applying the actions of each agent after the time step is finished. For this reason the API of PettingZoo [Ter+21] is based around turn-based environments. Examples of turn-based environments include most board games and many video games. Most real world scenarios are simultaneous though, so we only consider environments with simultaneous actions.

3.6. MARL Tasks in Related Work

Baker et al. [Bak+20] create a team-based multi-agent environment with one team of agents trying to find and “catch” the agents of the other team. Team et al. [Tea+21] create a generic physics-based world and a system to automatically generate cooperative and adversarial tasks like *Hide and Seek* and *Capture the Flag* within this world then train agents to be able to solve multiple tasks in this world. They then show their agents generalize to be able to solve unseen tasks.

Liu and Tan [LT20] test their attention-based architecture on three environments: (1) A catching game in a discrete 2D world, where multiple paddles moving in one dimension try to catch multiple balls that fall down the top of the screen. (2) A spreading game, where N agents try to cover N landmarks. The world is continuous, but the action space is discrete. This game is similar to the spreading game defined in [Low+17]. (3) StarCraft micromanagement: Two groups of units in the StarCraft game battle each other, each unit controlled by an agent.

Terry et al. [Ter+21] create an infrastructure for multi-task environments similar to the OpenAI Gym and add standardized wrappers for Atari multiplayer games, as well as a reimplementations of the MAgent environments [Zhe+18], the Multi-Particle-Environments [Low+17] and the SISL environments [GEK17].

Hüttenrauch, Šošić, and Neumann [HŠN19] test mean aggregation on two environments: (1) A task where all agents need to meet up in one point (rendezvous task). (2) A task where the agents try to catch one or more evaders (pursuit task). We compare our method on both of these tasks.

Gebhardt, Hüttenrauch, and Neumann [GHN19] define a set of environments based on Kilobots, simple small round robots in a virtual environment together with boxes. The Kilobots are tasked with moving the boxes around or arranging them based on their color. We compare our method on the box assembly as well as the box clustering task.

Chapter 4.

Scalable Information Aggregation for Deep MARL Policies

We introduce a policy architecture for deep multi-agent reinforcement learning that projects observations from one or more different kinds of observables into samples of latent spaces, then aggregates them into a single latent value. This makes the architecture scale to any number as well as a varying number of observables. We use the same encode-aggregate-concatenate structure as [HŠN19] and a similar method of attentive aggregation as [Bak+20].

4.1. Policy Architecture

While the policy architecture and weights are shared for all agents, the inputs are dependent on the current agent a_k , thus leading to different outputs for the action and value function for each agent.

Depending on the task, the inputs differ. In general, we always have the proprioceptive observations (how the agent a_k sees itself), and the observations from the n neighboring agents a_1, \dots, a_n . The observations from each neighbor have the same shape, but differ from the self-observations since they can include other information like the distance between a_k and a_i , and they may be missing some information that may not be known to a_k . In addition, there can be additional sets of observations, for example for objects or scripted agents in the environment. The architecture can handle any amount of observation sets.

Since we aggregate each set of observables of the same kind and shape into one aggregate, we call each of the sets of observations an “aggregation group” G :

$$G = \{g_1, \dots, g_n\}$$

G is a single aggregation group with each of the n aggregatables in the group called g_i . There can also be multiple aggregation groups that we call $(G^{(a)}, G^{(b)})$, each with their own number of observables. Since the g_i are later aggregated with a permutation-invariant aggregation function, their order does not matter.

First, we collect the observations for each instance $1 < i < n$ of each aggregation group G . These observations can be relative to the current agent. Each observation in the aggregation group G is thus a function of the agent a_k and the instance g_i . We call this observation $o_{k \rightarrow g_i}$, since agent k is observing the i th observable in group G :

$$o_{k \rightarrow g_i} = \text{observe}(a_k, g_i)$$

After collecting the observations for aggregation group g , the observations are each encoded separately into a latent embedding space:

$$e_{k \rightarrow g_i} = \text{enc}(o_{k \rightarrow g_i})$$

In general, the encoder $\text{enc}()$ could be an arbitrary differentiable function that maps the observation into a latent space. In our case, $\text{enc}()$ is a dense neural network with zero or more hidden layers, with shared weights across the observables in an observable group.

After being encoded, we interpret each instance of an aggregation group as one sample of a latent space that represents some form of the information of the aggregation group that is relevant to the agent. These samples are then aggregated using one of the aggregation methods described below to get a single latent space value for each aggregation group:

$$e_{k \rightarrow G} = \text{agg}_{i=0}^n(e_{k \rightarrow g_i})$$

We then concatenate all the latent spaces as well as the proprioceptive observations p_k of the agent k to get a single encoded value e_k :

$$e_k = (p_k || e_{k \rightarrow G^{(a)}} || e_{k \rightarrow G^{(b)}} || \dots)$$

This value is then passed through a decoder that consists of one or more dense layers. Finally, the decoded value is transformed to the dimensionality of the action space or to 1 to get the output of the value function. While we share the architecture for the policy and value function, we use a separate copy of the compute graph for the value function, so the weights and training are completely independent.

Figure 4.1 shows a schematic of the general model architecture described above.

The tasks we consider have a continuous action space. We use a diagonalized Gaussian distribution where the mean μ of each action is output by the neural network while the variance of each action is a free-standing learnable variable only passed through exp or softplus to ensure positivity.

4.2. Mean/Max/Softmax Aggregation

Each sample in the latent space is weighted by a function $\text{weigh}()$ and aggregated using an aggregation operator \oplus :

$$e_{k \rightarrow G} = \bigoplus_{i=1}^n \text{weigh}(e_{k \rightarrow g_i})$$

For mean aggregation, the weighing function multiplies by $\frac{1}{n}$ and the aggregation operator is the sum:

$$e_{k \rightarrow G} = \sum_{i=1}^n \frac{1}{n} \cdot (e_{k \rightarrow g_i})$$

For max aggregation, the weight is 1 and the aggregation operator takes the largest value:

$$e_{k \rightarrow G} = \max_{i=1}^n e_{k \rightarrow g_i}$$

For softmax aggregation, the weight is based on the softmax function and the aggregation operator is the sum:

$$e_{k \rightarrow G} = \sum_{i=1}^n \left(\frac{\exp(e_{k \rightarrow g_i})}{\sum_{j=1}^n \exp(e_{k \rightarrow g_j})} \right) e_{k \rightarrow g_i}$$

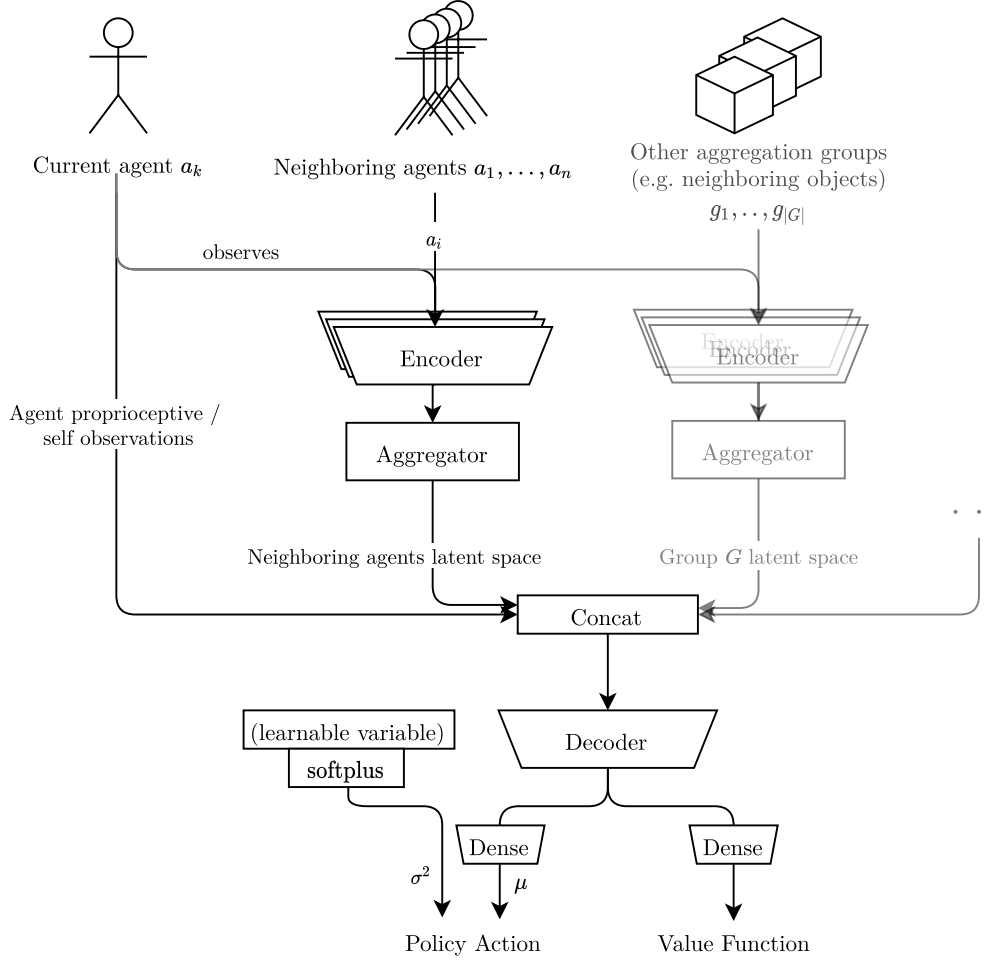


Figure 4.1.: A schematic of our general model architecture for deep MARL policies with scalable information aggregation. The observation inputs consist of the agent itself and multiple aggregatable observation groups. The observations from the aggregation groups are each passed through an encoder and aggregated with an aggregator. Afterwards all the aggregated observations are concatenated and decoded to get the policy and value function.

4.3. Bayesian Aggregation

We use a separate latent space and thus a separate observation model $p(z)$ for each aggregation group.

To make the Bayesian aggregation as described in Section 2.5.4 work in our policy network, we need an estimation describing the Gaussian prior (μ_{z_0} and $\sigma_{z_0}^2$) as well as the observed means r_i and variances $\sigma_{r_i}^2$.

The mean and variance of the Gaussian prior are learned as free-standing variables using the backpropagation during training, independent of the inputs of the neural network. The prior variance as well as the variance of the observations are rectified to enforce positivity using softplus.

To get the observed elements r_i , we use the two encoder networks enc_r and enc_σ that consist of a set of dense layers:

$$r_i = \text{enc}_r(o_{k \rightarrow g_i}), \quad \sigma_{r_i} = \text{enc}_\sigma(o_{k \rightarrow g_i})$$

enc_r and enc_σ are either separate dense neural networks or a single dense neural network with the output having two scalars per feature (one for the mean, one for the variance). Finally, we retrieve the value of $e_{k \rightarrow G}$ from the aggregated latent variable, using either just the mean of z :

$$e_{k \rightarrow G} = \mu_z$$

or by concatenating the mean and the variance:

$$e_{k \rightarrow G} = (\mu_z, \sigma_z^2).$$

The mean and variance are calculated from the conditioned Gaussian based on the encoded observations as described in Section 2.5.4:

$$\sigma_z^2 = \frac{1}{\frac{1}{\sigma_{z_0}^2} + \sum_{i=1}^n \frac{1}{\text{enc}_\sigma(o_{k \rightarrow g_i})^2}}$$

$$\mu_z = \mu_{z_0} + \sigma_z^2 \cdot \sum_{i=1}^n \frac{(\text{enc}_r(o_{k \rightarrow g_i}) - \mu_{z_0})}{\text{enc}_\sigma(o_{k \rightarrow g_i})^2}$$

A graphical overview of this method is shown in Figure 4.2.

4.4. Attentive Aggregation

When using residual self-attentive aggregation, we define the aggregated feature as

$$e_{k \rightarrow g} = \frac{1}{n} \sum_{i=1}^n \text{resatt}(\text{enc}(o_i))$$

with enc being a dense neural network with zero or more hidden layers, and

$$\text{resatt}(o_i) = o_i + \text{dense}(\text{MHA}(o_i, o_i, o_i)).$$

Residual means that the input is added to the output of the attention mechanism, so the attention mechanism only has to learn a *residual* value that modifies the input features.

The $\text{MHA}()$ module is the multi-head attention module from [Vas+17, sec. 3.2.2] as described in Section 2.5.5.

This method of attentive aggregation is similar to the method successfully used by Baker et al. [Bak+20]. An overview over the model architecture used in [Bak+20] can be seen in Figure 4.3.

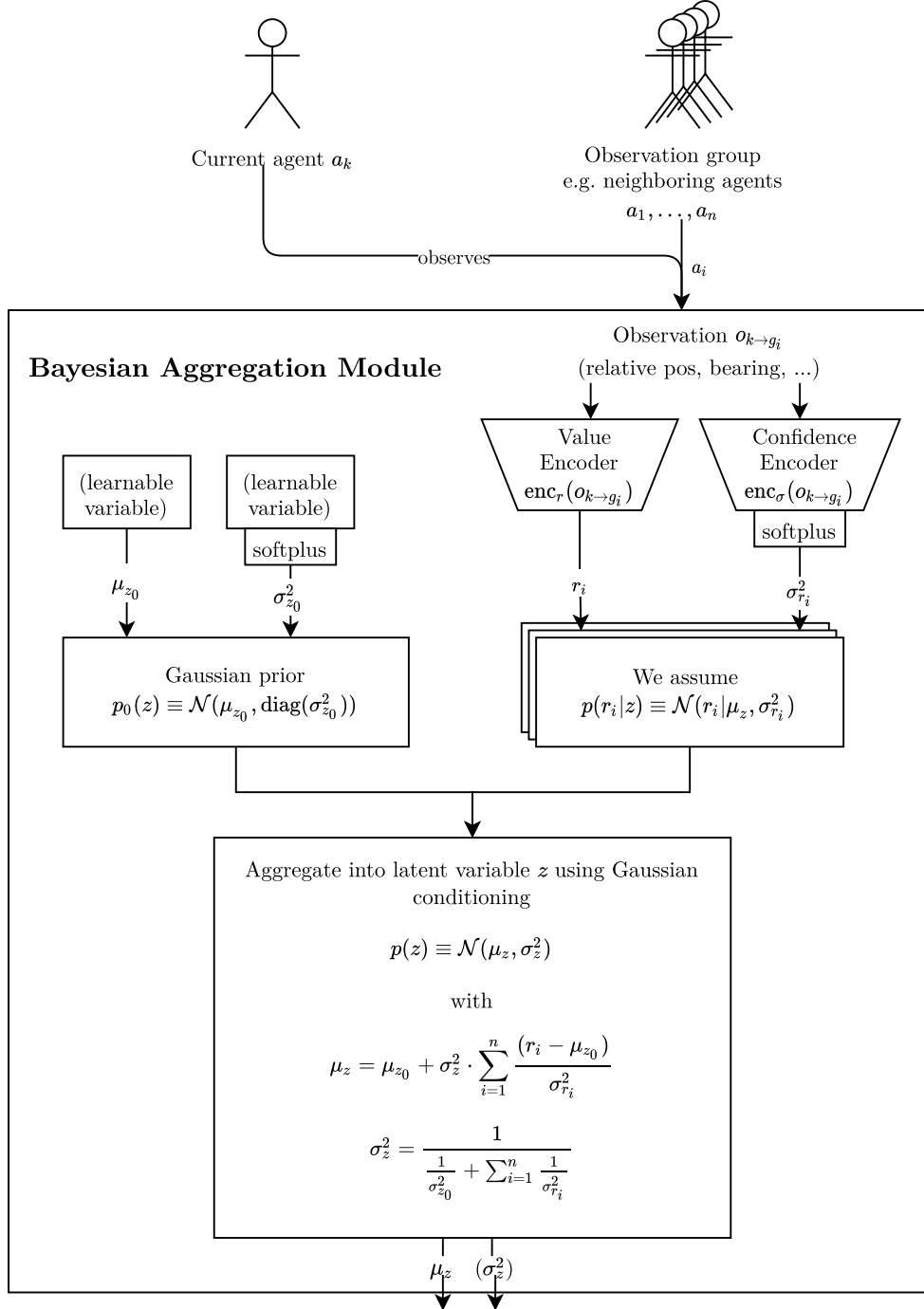


Figure 4.2.: Bayesian Aggregation in graphical form. The observations from the neighboring agents are encoded with the value encoder and the confidence encoder. They are then used to condition the Gaussian prior estimate of the latent variable z to get the final mean and variance of the a-posteriori estimate. The mean and optionally the variance estimate of z are concatenated with the latent spaces from the other aggregatables and passed to the decoder as shown in Figure 4.1.

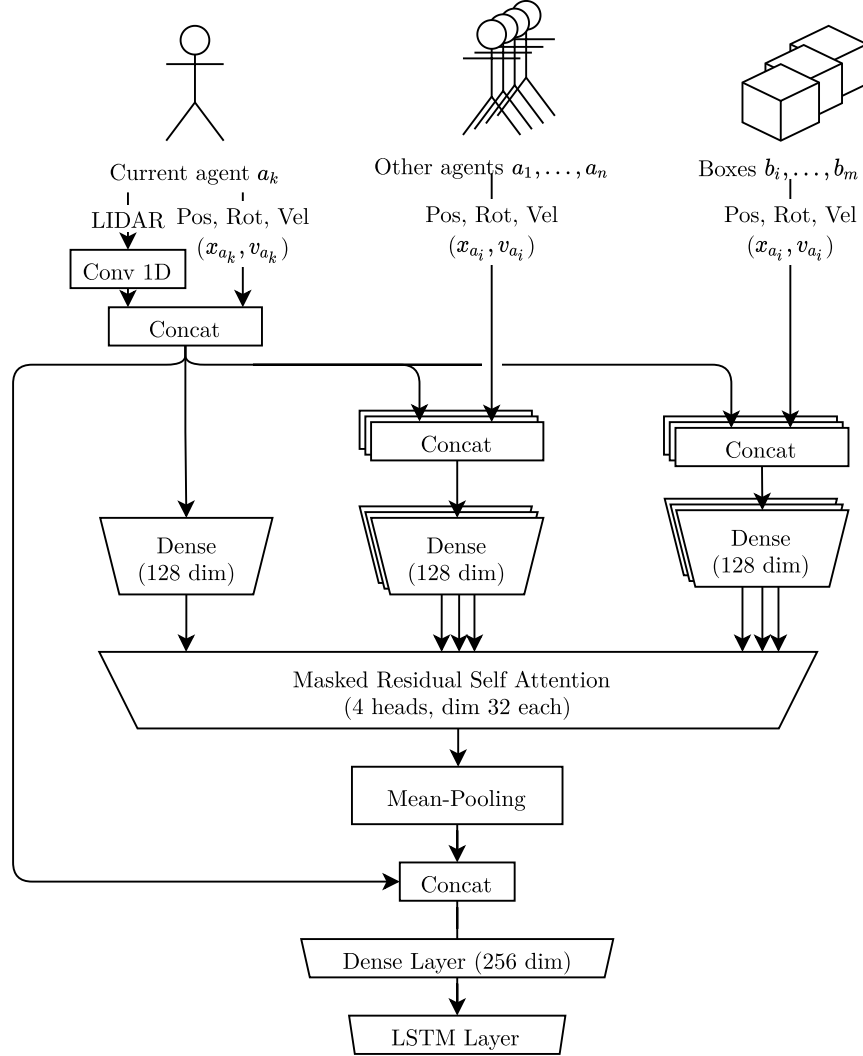


Figure 4.3.: A schematic of the model architecture used by OpenAI [Bak+20] using masked residual self-attention. It is similar to our architecture (Figure 4.1) except for the LIDAR in the self-observations as well as the LSTM layer at the end.

Chapter 5.

Experiments

5.1. Experimental Setup

All of our experiments were run using stable-baselines3 (sb3) [Raf+19]. *OpenAI Baselines* is a deep reinforcement learning framework that provides implementations of various reinforcement learning algorithms for TensorFlow. Stable-baselines is an extension of the original code base with better documentation and more flexibility for custom architectures and environments. Sb3 is the continuation of the stable-baselines project, rewritten using PyTorch. We rely on sb3 since the performance of PPO is known to depend on a number of implementation details [Eng+19] and the sb3 implementation of PPO-Clip has been shown to perform as well as the original OpenAI implementation [Raf+20].

We extended sb3 for the multi-agent use case by adapting the vectorized environment implementation to support multiple agents in a single environment, as well as adapting the training and evaluation functionality to correctly handle the centralized-learning decentralized-execution method. We also implement Trust Region Layers [Ott+20] as a new training algorithm for sb3 in order to be able to directly compare PPO and TRL. When mentioned in the results section, we optimized the hyperparameters using Optuna [Aki+19].

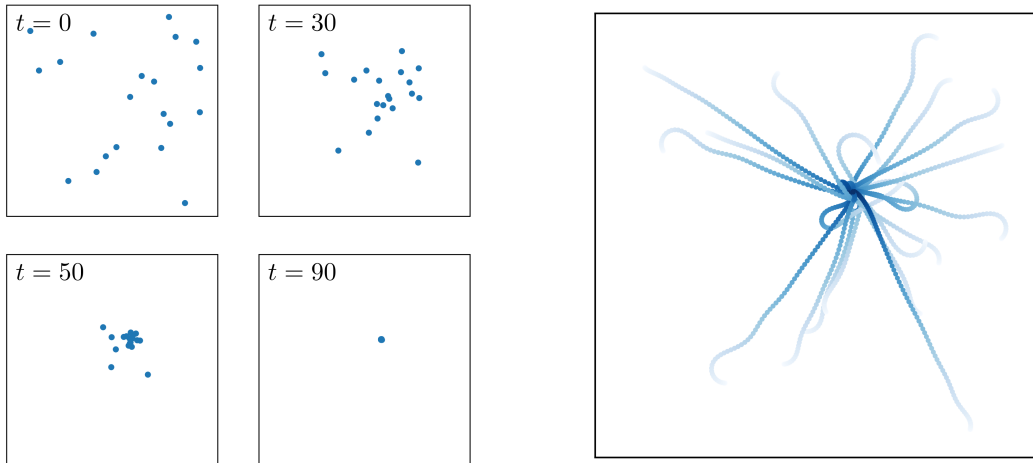


Figure 5.1.: Visualization of one successful episode of the rendezvous task (from [HŠN19])

5.2. Considered Tasks

To evaluate the different aggregation methods we need simulated environments where multiple agents can cooperatively solve a task. Since most of the commonly used tasks used to benchmark reinforcement learning algorithms are designed for a single agent, we use custom-built environments.

The following shows which specific tasks we consider. We start with a simple task and continue with progressively harder tasks that have more complex information in the environment that the policy needs to process.

5.2.1. Rendezvous Task

In the rendezvous task, a set of n agents try to meet up in a single point. An example is shown in Figure 5.1. Our implementation of the rendezvous task is modeled after [HŠN19].

The agents are modeled as infinitesimal dots without collisions. They use double-integrator unicycle dynamics [EH01], so the action outputs are the acceleration of the linear velocity (\dot{v}) and the angular velocity ($\dot{\omega}$). The agents move around in a square world that either has walls at the edge or is on a torus (the position is calculated modulo 100).

The observation space of agent a_i comprises the following information:

1. The own linear velocity v_{a_i}
2. The own angular velocity ω_{a_i}
3. The ratio of currently visible agents
4. If the observation space is not a torus: the distance and angle to the nearest wall
5. For each neighboring agent:
 1. The distance between the current agent and the neighbor
 2. The cos and sin of the relative bearing (the angle between the direction the agent is going and the position of the neighbor)
 3. The cos and sin of the neighbor's bearing to us
 4. The velocity of the neighbor relative to the agent's velocity
 5. The ratio of currently visible agents for the neighbor

The reward that every agent gets is the mean of all the pairwise agent distances:

$$r = \frac{1}{n} \sum_{i=0}^n \sum_{j=i+1}^n \|p_i - p_j\|$$

The episode ends after 1024 time steps.

5.2.2. Single-Evader Pursuit Task

In the pursuit task, multiple pursuers try to catch an evader. The evader agent has a higher speed than the pursuers, so the pursuers need to cooperate to be able to catch it. The world is a two-dimensional torus (the position (x, y) are floating point numbers modulo 100). If the world was infinite, the evader could run away in one direction forever, and if it had walls, the pursuers could easily corner the evader. The pursuit tasks are modeled after [HŠN19]. An example episode of the pursuit task is in Figure 5.2.

The agents are modeled with single-integrator unicycle dynamics. The action outputs are the linear velocity (v) and the angular velocity (ω).

The evader is not part of the training and uses a hard-coded algorithm based on maximizing Voronoi-regions as described by Zhou et al. [Zho+16]. It is thus part of the environment and not an “agent” as seen from the perspective of the training algorithm.

The observation space for the single-evader pursuit task comprises the following information:

1. The ratio of agents that are visible

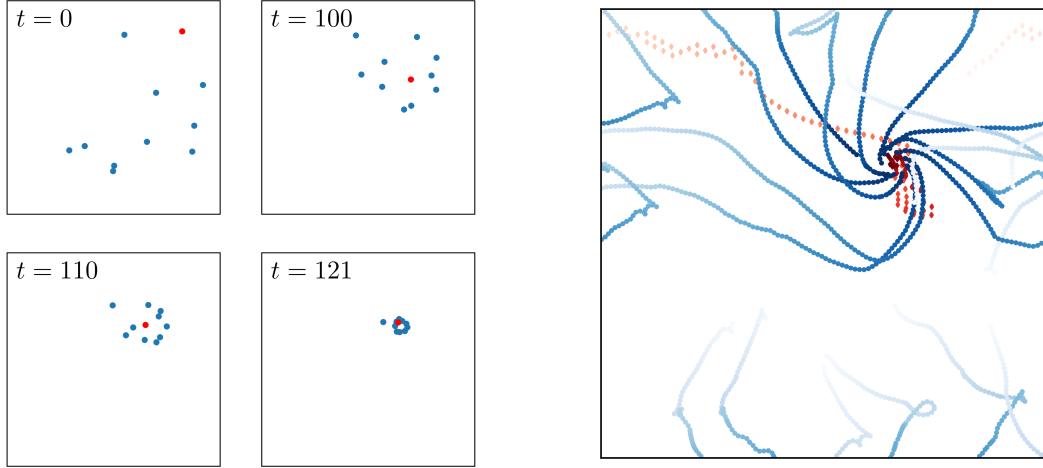


Figure 5.2.: Visualization of one successful episode of the pursuit task (from [HŠN19]). The pursuers are in blue, the evader is in red.

2. For each neighboring agent:
 1. The distance between the current agent and the neighbor
 2. The cos and sin of the relative bearing (the angle between the direction the agent is going and the position of the neighbor)
 3. The cos and sin of the neighbor's bearing to us
3. For the evader:
 1. The distance between the current agent and the evader
 2. The cos and sin of the relative bearing

The reward function of all the agents is the minimum distance of any agent to the evader:

$$r = \min_{i=0}^n ||p_{a_i} - p_e||$$

The episode ends once the evader is caught or 1024 time steps have passed. The evader is declared as caught if the distance is minimum distance between an agent and the evader is less than 1% of the world width.

The observation of the evader is handled as a single observable in a second aggregation group to be able to scale the policy architecture to multiple evaders without modifications.

5.2.3. Multi-Evader Pursuit

The multi-evader pursuit task is the same as the normal pursuit task, except there are multiple hard-coded evaders.

The reward here is different, since defining the reward as for the single-evader task is not obvious. The reward is $+1$ whenever an evader is caught, 0 otherwise. An evader is declared caught if the distance to the nearest pursuer is less than $\frac{2}{100}$ of the world width. Contrary to the single-evader task, the episode does not end when an evader is caught and instead always runs for 1024 timesteps.

5.2.4. Box Assembly Task

In the box assembly task, the agents are modeled similar to Kilobots, as described by Gebhardt, Hüttenrauch, and Neumann [GHN19]. The agents are two-dimensional circles with collision.

For simplicity, we again use single-integrator unicycle dynamics with the linear velocity v and the angular velocity ω as the action outputs instead of the specific dynamics of real Kilobots.

The world is a square and contains a few two-dimensional boxes (squares) as obstacles. For the box assembly task, the goal is to get all the boxes as close together as possible. Since the agents are much smaller than the boxes, moving a box is hard for a single agent. The reward is the derivation of the negative sum of the pairwise distances between the boxes. We run this task with four boxes and 10 agents.

An example of a successful episode of the task is in Figure 5.3.

The observation space for the box assembly task contains the following information:

1. The absolute position (x, y) of the current agent
2. The sin, cos of the absolute rotation of the current agent
3. For every neighboring agent:
 1. The distance between the current agent and the neighbor
 2. The sin and cos of the bearing angle to the neighbor
 3. The sin and cos of the neighbor's rotation relative to our rotation
4. For every neighboring box:
 1. The distance between the current agent and the box
 2. The sin and cos of the bearing angle to the box
 3. The sin and cos of the box's rotation relative to our rotation

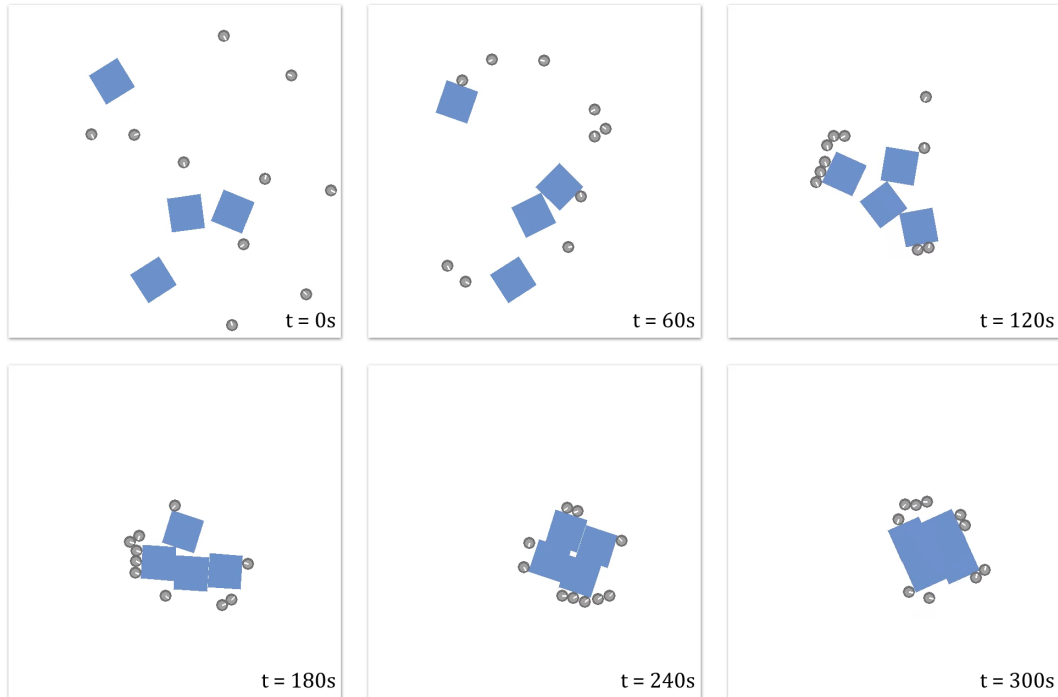


Figure 5.3.: Example successful episode of the box assembly task (from [GHN19])

5.2.5. Box Clustering Task

The task setup for box clustering is the same as for the box assembly task, except that each box is assigned a color. The goal is to move the boxes into an arrangement such that boxes of the same color are as close together as possible, while boxes of different colors are far away. Each color of box has an explicit target in one corner of the world. The reward is calculated by taking the sum of negative distances between each box and its cluster target, then taking the difference of this after the step and before the step.

The observation space is the same as the observation space of the box assembly task, except that each cluster of objects is aggregated into a separate aggregation space, and that there is an additional one-hot encoded cluster ID in the object observation space. Figure 5.4 shows an example successful episode of the clustering task with two four boxes split into two clusters.

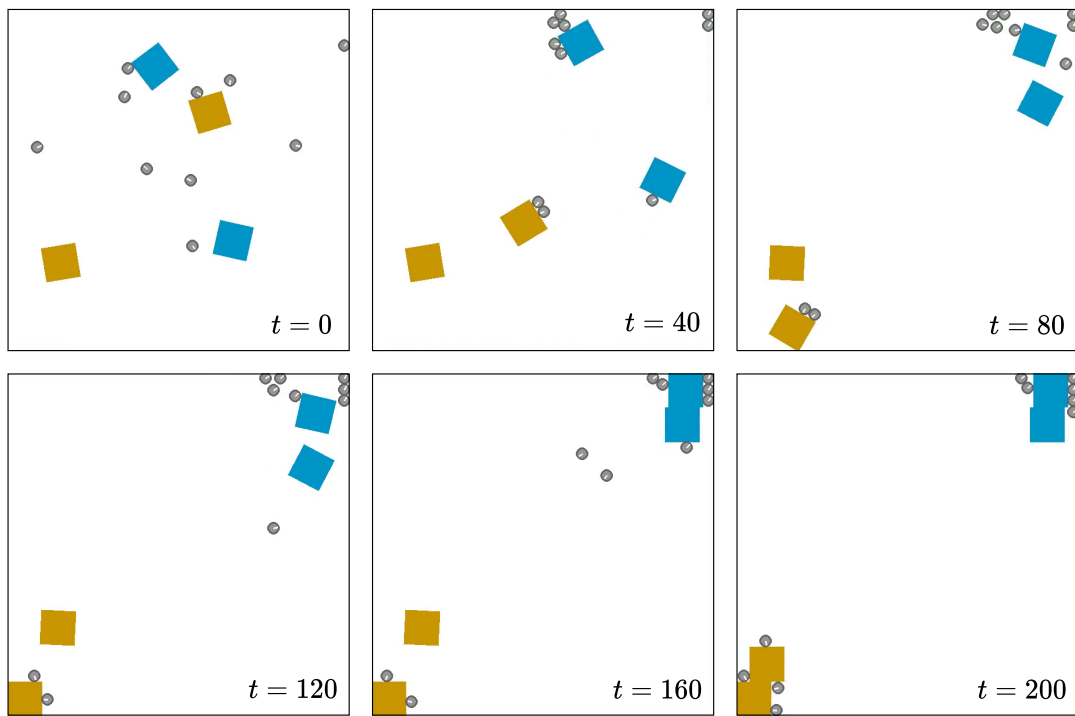


Figure 5.4.: Example episode of the clustering task with two clusters.

Chapter 6.

Results

In this section, we present the results of our experiments. We compare (1) the training performance of the different aggregation methods on various tasks, (2) aggregating different observation sets into a single vs multiple latent spaces (3) various variants of Bayesian aggregation, and (4) the training algorithms PPO and TRL.

In order to compare the performance of different variants, we consider the average reward per episode over the agent steps used during training. This way we can evaluate and compare the overall best results as well as the sample efficiency.

Each policy gradient training step consists of multiple gradient descent steps with a set of batches of trajectory rollouts generated using the previous policy. We evaluate the performance after each training step on a separate set of evaluation environments.

Unless mentioned otherwise, the following setup is used:

- The result plots show the median reward of the n runs at a specific training step, the error band is the 25th and 75th percentile of runs
- The activation function used after each layer is LeakyReLU
- For the Bayesian aggregation:
 - We only use the mean value μ_z as an output, not μ_z and σ_z^2
 - We use a single shared encoder for the value and confidence estimates
 - The a-priori estimate μ_{z_0} is learnable separately for each feature dimension
 - The variance of the a-priori estimate as well as the encoded estimates are rectified using softplus.

- Multiple aggregatable groups are aggregated into separate latent spaces
- The parameters of the policy and value function are not shared
- The training algorithm used is PPO

Due to time and resource constraints, we only use individually hyper-parameter optimized architectures on the multi-evader pursuit and rendezvous tasks, and use the same architecture for all aggregation methods on the other tasks.

6.1. Aggregation Method Results

We compare the performance of the different aggregation methods (Bayesian aggregation, mean aggregation, attentive aggregation) on multiple tasks.

We use the notation **64-agg-64** to describe the layer sizes of the neural network: The numbers before **agg** are the sequential layer sizes of the dense layers of the encoders of each aggregation group. The numbers after **agg** are the layer sizes in the decoder after the concatenation of the proprioceptive observations with the aggregated observations (compare Figure 4.1).

We start with a very simple task (rendezvous), then look at progressively more difficult tasks with more complex information that needs to be aggregated.

6.1.1. Rendezvous Task

The rendezvous task is very simple and thus gives a good baseline to show whether our model works in general. It only has one simple aggregation group (the agents).

Figure 6.1 shows a comparison between mean and Bayesian aggregation on the rendezvous task with twenty agents in a two-dimensional square world with walls. The medium architecture is the one optimized on the pursuit task (**120-60-agg-160**). The simple architecture is the one used in [HŠN19] (**64-agg-64**). The optimized architecture was optimized on the rendezvous task with Bayesian aggregation directly: **146-120-agg-19-177-162**. All architectures and aggregation methods successfully solve the task. The aggregation method does not make any difference in the performance, but both the simple and the medium architecture have a “drop” in training speed at around 2 million steps, while the optimized architecture smoothly learns the problem. The logarithmic scale graph to the right shows that while the simple and optimized architecture both reach the same final score, the medium architecture never reaches the same score. This might be because both the simple and the optimized architecture have a bottleneck layer after the aggregation, forcing the neural network to simplify the strategy.

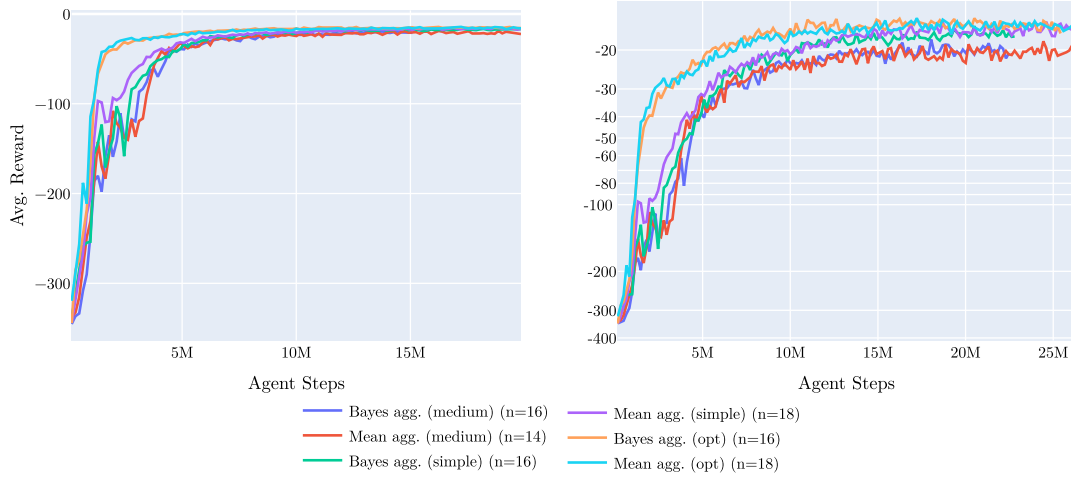


Figure 6.1.: Results on the rendezvous task. *Opt* is the hyper-parameter optimized architecture. The results barely differ between the mean and Bayesian aggregation, but the size of the policy architecture makes a difference. In the logarithmic view on the right it can be seen that the medium architecture does not reach the same final performance as the other two architectures.

We show a random exemplary episode of each a policy with median performance for the medium and optimized architectures in Figure 6.2. For the medium architecture case with the Bayesian aggregation, we notice that the consensus location does not stay the same and instead drifts even after all agents are at the same spot. In addition, sometimes there is one or more stragglers that take significantly longer to get to the consensus location. For the medium architecture with the mean aggregation, the agents do not converge on a single point, but instead continue to circle a small area. For the optimized architecture, the consensus is reached quickly (after 100 of 1000 time steps) and stays at the same spot.

6.1.2. Single-Evader Pursuit Task

The pursuit task is a slightly more complex task with the single evader being in a second aggregation group with one observable. The results here show that the learned policies can handle a moving object in the world and that the architecture works with multiple observation groups.

Figure 6.3 shows the results on the single-evader pursuit task with 10 pursuers and one evader. The neural network architecture is fixed at 120-60-agg-160 for all methods. All methods learn the task quickly with the same best performance, with the mean aggregation achieving the maximum performance slightly faster. This shows that the task is simpler than the multi-evader pursuit task, even though the policy

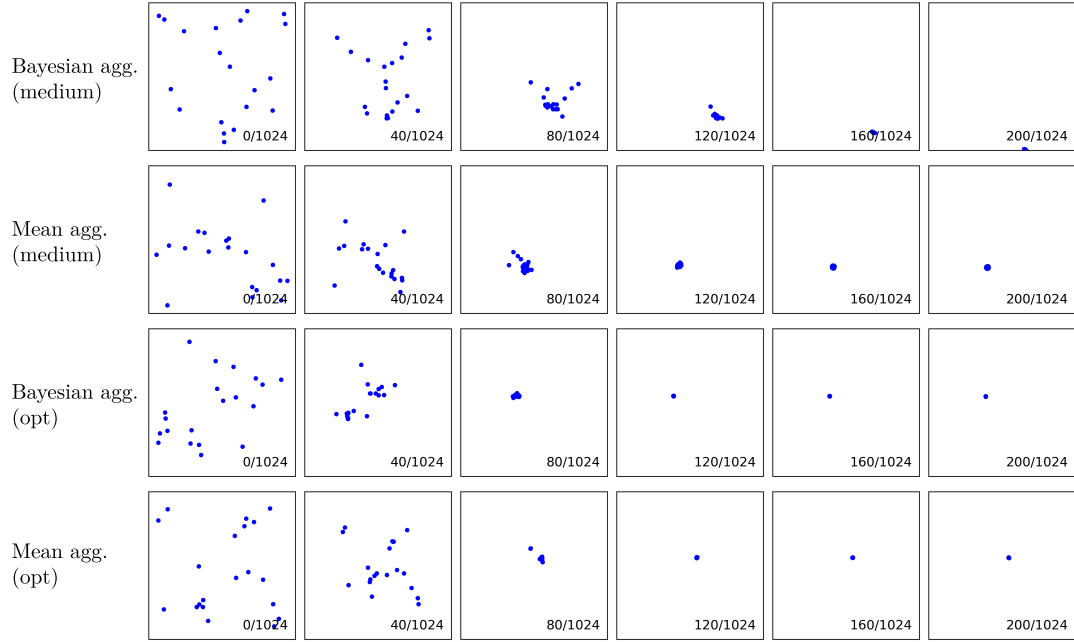


Figure 6.2.: View of a random episode of the rendezvous task with the medium and optimized architectures. *Bayesian agg. (medium)* tends to drift and have stragglers even after a consensus seems to have been reached. *Mean agg. (medium)* does not converge to an exact point, the agents continue to circle a small area. Both aggregations with the *opt* architecture converge quickly to a single point and stay in one place.

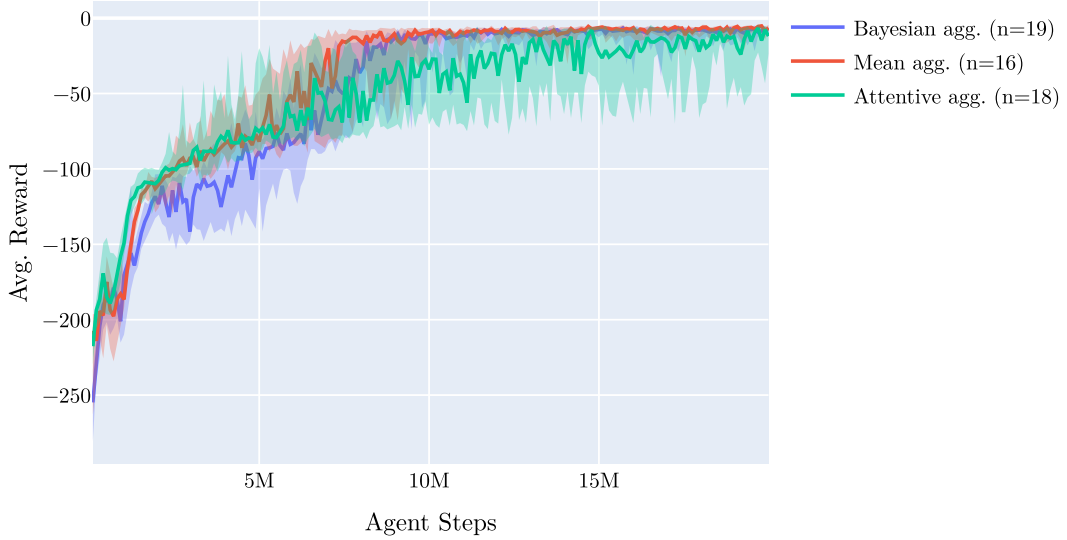


Figure 6.3.: Results on Single Pursuit task. The performance is similar for all methods, with the mean aggregation achieving the best performance slightly faster.

architecture is the same. This is both due to the fact that there are fewer evaders and that the reward is more sparse (minimum-distance for single-evader vs count-catches for multi-evader).

6.1.3. Multi-Evader Pursuit Task

The multi-evader pursuit task adds complexity by having multiple evaders that need to be caught. In addition, the reward is sparser due to the fact that each catch only gives a binary reward signal. The results here show whether our policy is able to process multiple aggregation groups with multiple moving entities. Here, we consider the multi-evader pursuit task with 50 pursuers and 5 evaders on a torus. Figure 6.4 shows the results of the multi-evader pursuit task with different aggregation methods with the same architecture used in [HŠN19] to be able to directly compare the results. The architecture is 64-agg-64 with the tanh activation function. With this architecture, the Bayesian aggregation performs best.

Figure 6.5 shows the results with neural network architectures that were separately hyper-parameter optimized for each aggregation method except max aggregation. The optimized architecture for the mean aggregation is 174-226-97-agg-96, the same architecture is used for the max aggregation. The optimized architecture for the Bayesian aggregation is 120-60-agg-160. The optimized architecture for the attentive aggregation is 72-agg-132-200. All optimized architectures use the LeakyReLU activation functions. Note that the architecture optimized on the mean aggregation is the deepest with three hidden layers before the aggregation, while the optimized



Figure 6.4.: Results on the multi-evader pursuit task with the NN architecture adapted from [HŠN19]. The Bayesian aggregation performs best.

architecture on the attentive aggregation has multiple layers after the aggregation instead. With the hyper-parameter optimized architecture, the mean aggregation performs best. The results are still similar when using the same 120–60–agg–160 architecture for every aggregation method. The results of Figures 6.4 and 6.5 indicate that the Bayesian aggregation outperforms the mean aggregation when the neural network is limited in size, but has no advantage when the neural network is sufficiently large and deep. The neural network seems to be able to implicitly learn to transform and weigh the information from the different observables, compensating the advantage of the additional structure of relevancy / certainty that is given in the Bayesian aggregation.

Figure 6.6 shows the results of the top third of runs. The performance of the mean and Bayesian aggregation is similar, indicating that the best runs are similar, but that the Bayesian aggregation has more runs that fail to achieve the peak performance.

Figure 6.7 shows the first 500 steps of a random episode with the medium-performing policy and the best-performing policy of each of the aggregation methods. The mean aggregation tends to learn policies where first all the evaders are forced to go to one location, then they are circled and caught all at once. Since the evaders are algorithmically controlled and always go to the largest free Voronoi region, the evaders naturally go to the same location if they have enough space, so this strategy works well. The strategy of the Bayesian aggregation policy seems to be more flexible, where different groups of pursuers catch different groups of evaders at the same time instead of all concentrating on a single cluster of evaders. This strategy does not lead to better overall performance, however.

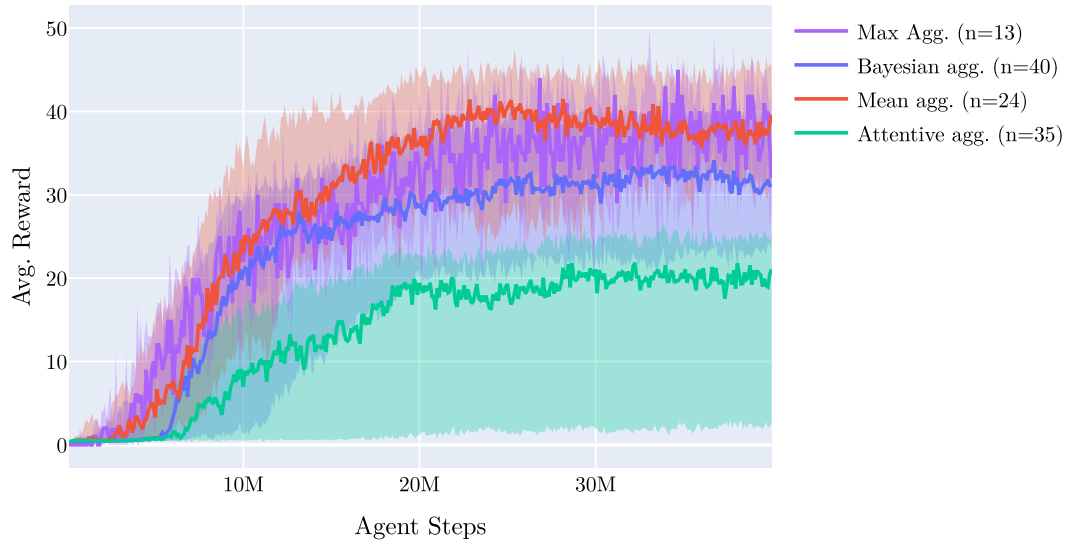


Figure 6.5.: Results on the multi-evader pursuit task with the NN architecture hyperparameter optimized for each aggregation method. The mean aggregation performs best.

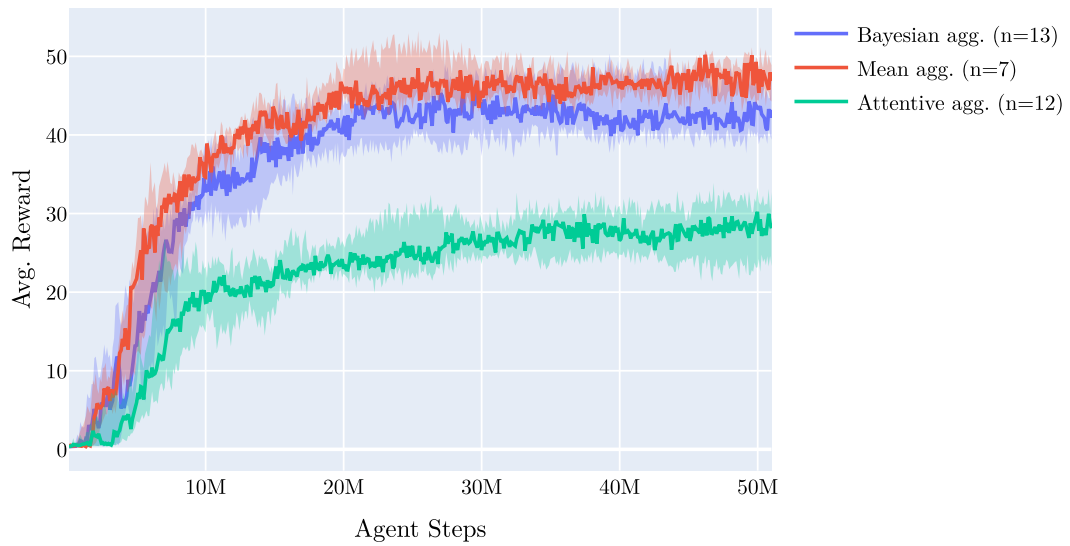


Figure 6.6.: Like Figure 6.5 but only the top 1/3 of runs. This shows that the peak performance of the mean and the Bayesian aggregation is similar.

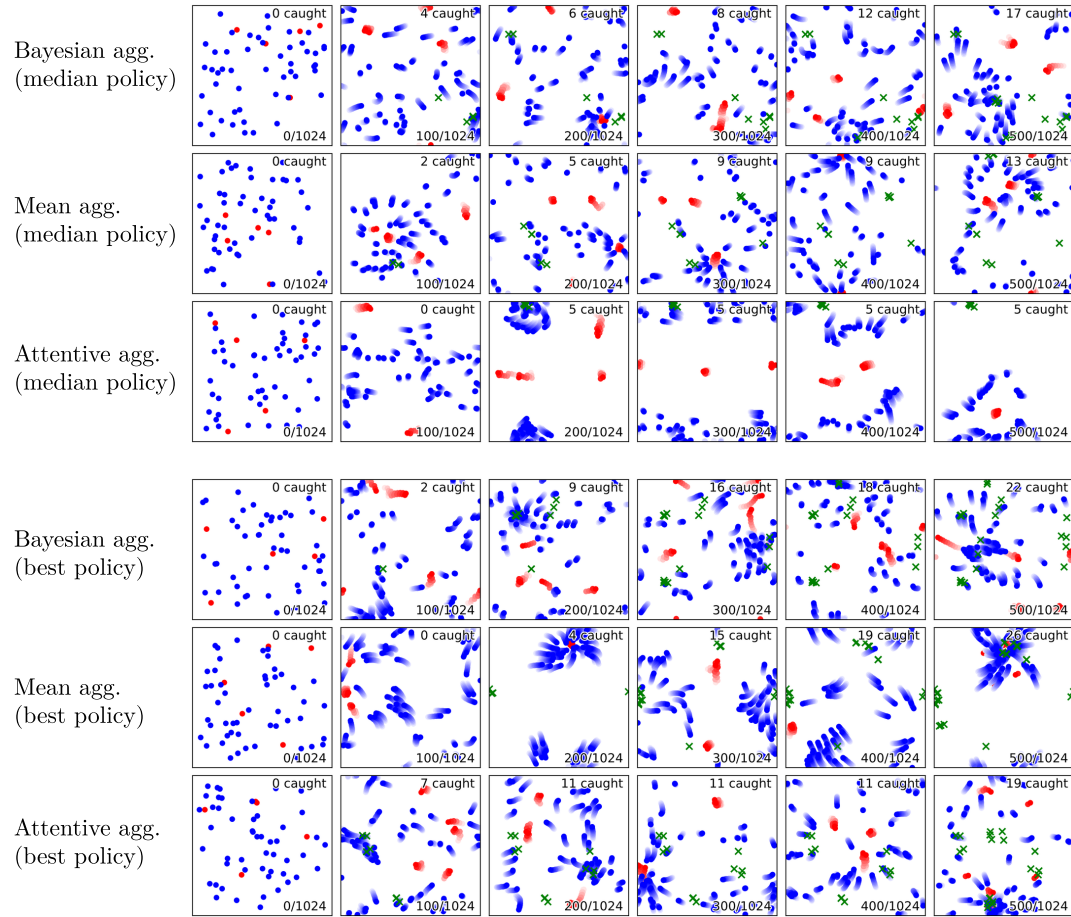


Figure 6.7.: Random episode with the median and best performing policy of each aggregation method. Evaders are respawned when they are caught. Locations where an evader is caught are marked in green, with the total catch count up to the current time step is in the top right. The best performing policies of the different aggregation methods seem to have different strategies.

6.1.4. Assembly Task

The assembly task is different from the previous tasks in that the agents have to manipulate the environment in order reach their goal. In order to solve the task, the agents need to show collaborative skills in deciding on the final destination of the assembled cluster of objects as well as succeed in moving the objects to that location. They also need to understand when the goal has been reached in order to not move the boxes apart again. Figure 6.8 shows the results on the assembly task with ten agents and four boxes. The three aggregation methods perform very similar, with the attentive aggregation learning the task slightly quicker.

6.1.5. Clustering Task With Two Clusters

The clustering task is an extension of the assembly task with a *third* aggregation group for the second cluster of objects. To successfully solve this task, the agents need to be able to distinguish between the two object types and split the work of moving them. Figure 6.9 shows the results on the clustering task with four boxes split into two clusters. The results are similar in all cases, with all policies solving the task for most runs. An example of a successful episode with mean aggregation is shown in Figure 5.4.

6.1.6. Clustering Task With Three Clusters

We can extend the clustering task to have any number of boxes in any number of clusters, making it more difficult. We find that when extending the task to three clusters, none of our trained policies manage to adequately solve the task. In most cases, the agents only move the boxes of one cluster, manage to mostly move those boxes to their target location, but then stay in the corner completely ignore the boxes in the other clusters. Only when increasing the number of trajectories per training step and the batch size to 10x the size used in the other experiments does a policy consider boxes in multiple clusters (see Appendix A). ?? shows an episode with some unsuccessful episodes and with the best policy trained with Bayesian aggregation. Even the best policy only considers two of the three clusters and ignores boxes in the third cluster. Since the task is not solved, we do not provide a comparison graph of the different aggregation methods.

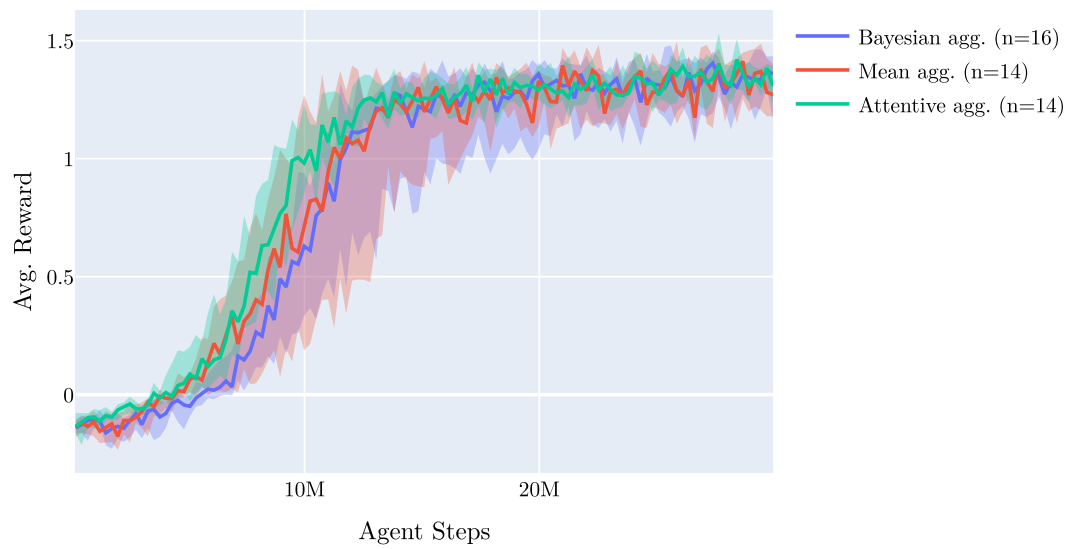


Figure 6.8.: Results on the assembly task.

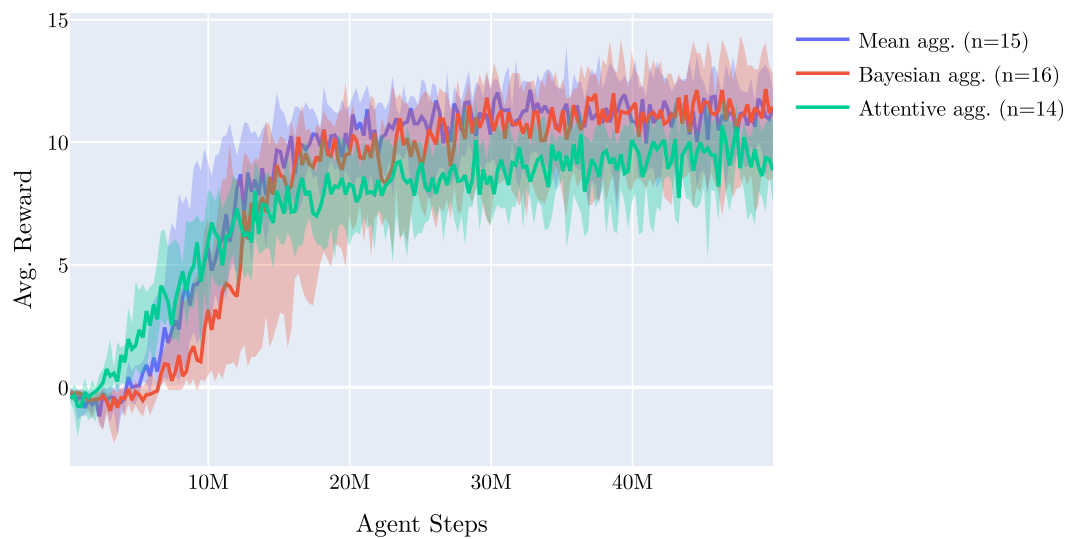


Figure 6.9.: Results on the clustering task with two clusters. The attentive aggregation starts learning slightly faster, but has a lower final result.

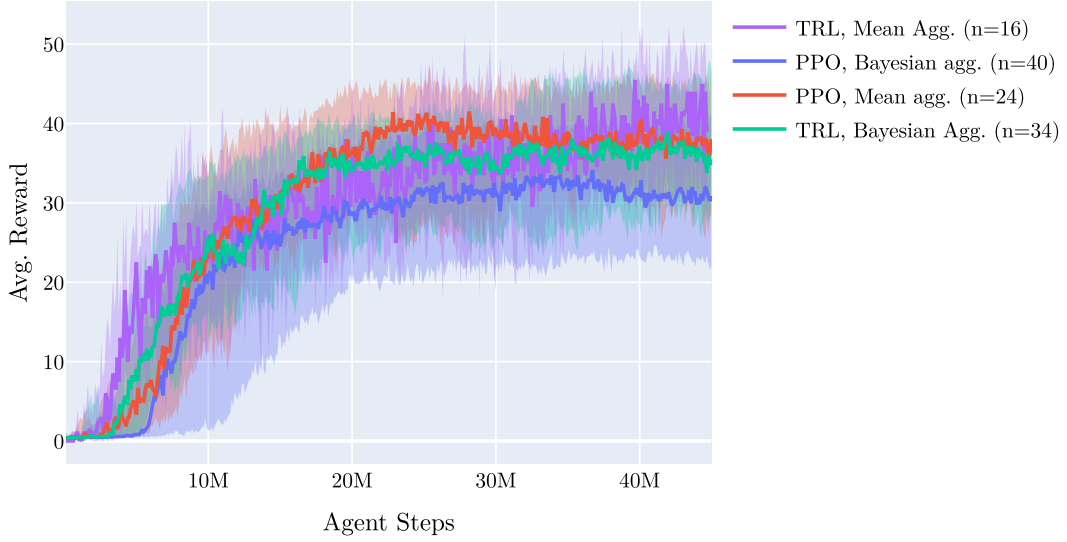


Figure 6.10.: TRL vs PPO learning algorithms on the multi-evader pursuit task. For Bayesian aggregation, the training is faster, and the end result better. For mean aggregation, the results are fairly similar.

6.2. Learning Algorithm Comparison (PPO vs PG-TRL)

In multi-agent fully cooperative tasks, the reward can be even more sparse than it is for single-agent tasks, since the same reward is used for all agents. One agent can thus not know whether its own actions or the actions of a different agent lead to a specific reward signal. Trust region layers policy gradient (TRL) (see Section 2.2.2) was published for single-agent reinforcement learning, but since TRL promises to improve the exploration over PPO and possibly find a better optimal solution, we apply it to our multi-agent experiments in a direct comparison to PPO.

In the following, we show some training graphs of the TRL training method compared to PPO. Figure 6.10 shows the learning algorithm comparison on the multi-evader pursuit task. The architectures are the ones hyper-parameter optimized on PPO on each of the aggregation methods. TRL seems to show significantly improved training performance for the Bayesian aggregation and similar performance for the mean aggregation.

Figure 6.11 shows the comparison on the single-evader pursuit task. Here, the performance of TRL is better than PPO on both the mean and the Bayesian aggregation.

Figure 6.12 shows the comparison on the assembly task. All aggregation methods perform the same, except the Bayesian aggregation with TRL, which is worse.

In summary, TRL seems to perform the same or better than PPO in most cases, with the Bayesian aggregation on the assembly task being an outlier. Since we did

Single-evader Pursuit TRL vs PPO

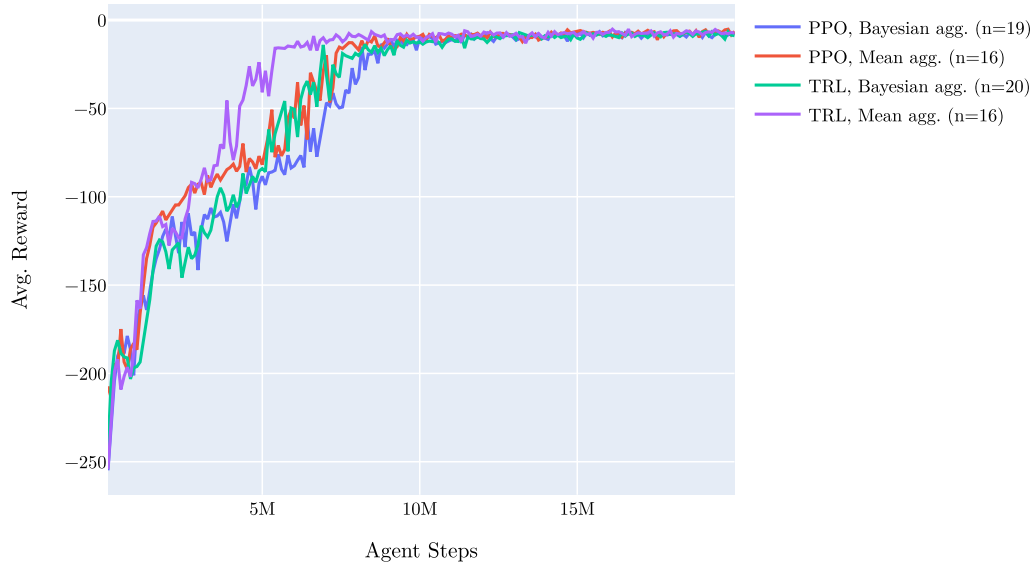


Figure 6.11.: TRL vs PPO learning algorithms on single-evader pursuit task. The performance of TRL is better for both the mean and the Bayesian aggregation.



Figure 6.12.: TRL vs PPO learning algorithms on the assembly task. The training performance is the same for all aggregation methods, except that the Bayesian aggregation performs worse with TRL.

not do extensive testing with different hyper-parameters for TRL, we can conclude that TRL may have good potential of performing better, in general, than PPO on multi-task environments with sparse rewards.

6.3. Bayesian Aggregation Variants

In the results comparing the different aggregation methods above, we always use the same setup for the Bayesian aggregation. The following shows the results for some variants of the Bayesian aggregation on selected tasks.

6.3.1. Single Space vs Separate Space Aggregation

For tasks where we have multiple aggregation groups, we can also aggregate all observables into the same latent space instead of separate ones. This means that instead of each aggregation space containing the information of those aggregatables, the single aggregation space must contain the information of all observables as it pertains to the current agent. The potential advantage is that the policy can share more parameters and thus be more sample efficient. It can also scale to a larger number of categories of objects (aggregation groups). Figure 6.13 shows a schematic comparison between the two methods.

In the other experiments we always use separate spaces. Figure 6.14 shows the results of aggregating into a single space vs. into separate spaces for the multi-evader pursuit and assembly tasks. In the case of the multi-evader pursuit task single-space aggregation means that both the neighboring pursuers and all the evaders are aggregated into one space. For the assembly task, both the neighboring agents and the boxes are aggregated into one space. The separate-space aggregation performs better in both tasks, with a higher margin in the multi-evader task.

In summary, while the single-space aggregation is promising by in theory being able to scale to more complex environment, it performs worse for our tasks. This might be due to the fact that learning different encoder networks to output information into the same latent space is hard. An alternative explanation might be that our experiments only have few aggregation groups, so the value of parameter sharing is low, or due to the lack of experiments with more different hyper-parameters.

6.3.2. Separate vs Common Encoder

As described in Section 4.3, we can either have a shared encoder to predict the mean and variance of each sample in each aggregation space by making the last layer of

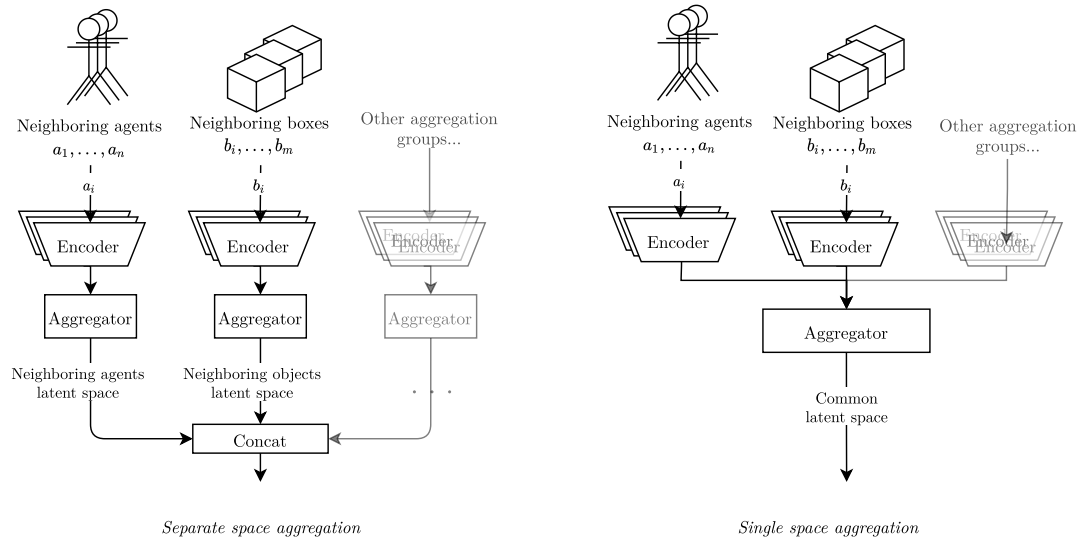


Figure 6.13.: Schematic comparison of separate space vs. single-space aggregation. In single-space aggregation the encoders are still separate, but the latent space is shared. The single-space aggregation shares more parameters and can scale better to more aggregation groups.

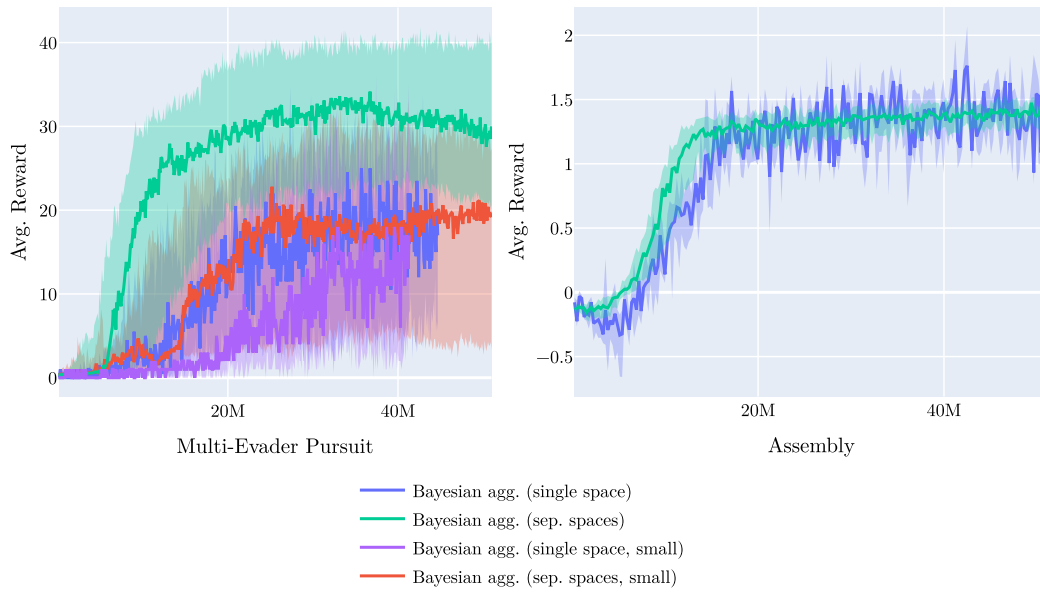


Figure 6.14.: Separate vs single-space aggregation on the multi-evader pursuit and assembly tasks. The single-space aggregation performs worse in both tasks.

the encoder have two outputs for each feature, or have two fully separate networks (enc_r and enc_σ). In our experiments, using one common encoder with two outputs generally performs better.

6.3.3. Using the Aggregated Variance or Only the Mean

In the other experiments with Bayesian aggregation, we only use the predicted mean of the Gaussian distribution of each latent space feature as an input to the decoder:

$$e_{k \rightarrow G} = \mu_z$$

Since the Bayesian aggregation also gives us a full a-posteriori Gaussian distribution, we also have an estimate of the variance for each feature in the latent space that is computed from the apriori variance conditioned on each seen latent space sample. We can feed this variance to the decoder by concatenating it with the mean predictions in the hope that the neural network is able to use this additional information:

$$e_{k \rightarrow G} = (\mu_z, \sigma_z^2)$$

The results of applying this method to the multi-evader pursuit are seen in Figure 6.15. The neural network architecture is the same as for the other experiments with Bayesian aggregation on multi-evader pursuit (120-60-agg-160). Including the variance in the decoder inputs decreases the performance.

The decreasing performance could be a result of the increased dimension of the decoder inputs. Adding the variance inputs doubles the number of values the decoder has to process and learn from. Since the structure of the encoded values and variances can not be known beforehand to the decoder, it has to learn to interpret more information than when receiving just the mean values. The added variance inputs should give the decoder the ability to understand the confidence of each of the value predictions and weigh them accordingly, but the added complexity seems to make it not worth it.

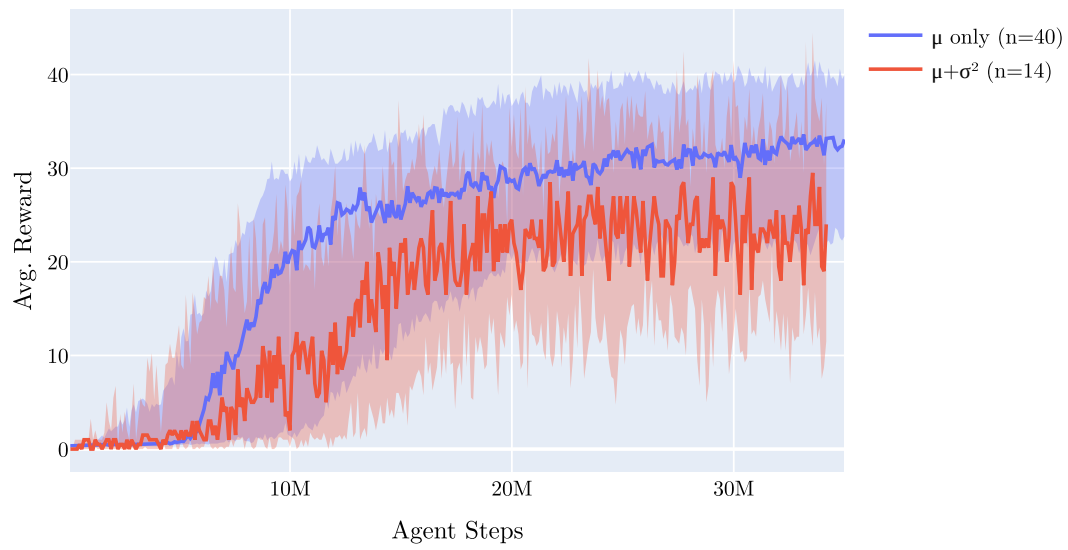


Figure 6.15.: Results of Bayesian aggregation on the multi-evader pursuit task, depending on whether the variance is also fed into the decoder or only the mean.

Chapter 7.

Conclusion and Future Work

7.1. Conclusion

We have made a comprehensive comparison between the performance of mean aggregation, Bayesian aggregation, and attentive aggregation to collect a varying number of observations on a set of different deep reinforcement learning tasks. We have observed that there is no clear advantage of one of the methods over the others, with the results differing strongly between the different tasks. In addition, the signal to noise ratio of the comparisons was fairly low, since other hyperparameters like the sizes of the neural networks or the training step size changed the results more than did the chosen aggregation method.

We have also shown the results of a few variants of the Bayesian aggregation and concluded that it performs best when (1) encoding the variance with the same encoder as the estimate instead of with separate encoders, (2) aggregating observables of different kinds into separate latent spaces instead of the same one and (3) not using the aggregated variance as an input to the decoder.

Finally, we have applied a new training method (trust region layers) to multi-agent reinforcement learning and compared it to the commonly used PPO. The results indicate that TRL is usually at least as good as PPO and in some cases outperforms it, indicating that it may be a good alternative to PPO for environments with sparse rewards, as is the case in many MARL tasks.

7.2. Future Work

There are many avenues for future work in this area.

All of our experiments used global visibility due to the noisy nature of limited local visibility making it harder to make any strong conclusions. Since the local visibility case also increases the uncertainty of each observation though, it might be a case where Bayesian aggregation performs better. Future work should include experiments that have a larger environment with observability limited by range or by obstacles.

Even though we show in our experiments that Bayesian aggregation into the same latent space is worse than aggregating into separate spaces, there might be potential with this and the other variants of Bayesian aggregation we introduced to be able to scale the number of aggregation groups more. This could be accomplished with more hyperparameter tuning or by introducing a two-stage encoder where the first encoder is separate by aggregation group and the second encoder is shared, then aggregating the output of the second encoder into the same space.

We also only considered tasks with implicit communication — the agents had to infer the intent of the other agents purely by their actions. There is related work that adds explicit communication between agents that is learned together with the policy. This is usually implemented as another action output that is written directly into the observation of the other agents instead of affecting the world. Explicit communication architectures may be able to handle some tasks better than those with implicit communication, but they are often only applicable to environments with exactly two agents. For more agents, the performance of explicit communication architectures may be affected by the aggregation methods used and thus might be a use case for Bayesian aggregation.

Bibliography

- [Aki+19] Takuya Akiba et al. “Optuna: A Next-generation Hyperparameter Optimization Framework”. In: *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. KDD ’19. Anchorage, AK, USA: Association for Computing Machinery, July 2019, pp. 2623–2631. ISBN: 9781450362016. DOI: 10.1145/3292500.3330701. URL: <https://doi.org/10.1145/3292500.3330701> (visited on 07/07/2021).
- [Bak+20] Bowen Baker et al. “Emergent Tool Use From Multi-Agent Autocurricula”. In: *arXiv:1909.07528 [cs, stat]* (Feb. 2020). URL: <http://arxiv.org/abs/1909.07528> (visited on 06/14/2021).
- [Bis06] Christopher Bishop. *Pattern Recognition and Machine Learning*. Information Science and Statistics. New York: Springer-Verlag, 2006. ISBN: 9780387310732. URL: <https://www.springer.com/gp/book/9780387310732> (visited on 06/14/2021).
- [Bür09] Axel Bürkle. “Collaborating Miniature Drones for Surveillance and Reconnaissance”. In: *Unmanned/Unattended Sensors and Sensor Networks VI*. Vol. 7480. International Society for Optics; Photonics, Sept. 2009, 74800H. DOI: 10.1117/12.830408. URL: <https://www.spiedigitallibrary.org/conference-proceedings-of-spie/7480/74800H/Collaborating-miniature-drones-for-surveillance-and-reconnaissance/10.1117/12.830408.short> (visited on 07/20/2021).
- [Can+21] Lorenzo Canese et al. “Multi-Agent Reinforcement Learning: A Review of Challenges and Applications”. In: *Applied Sciences* 11.11 (Jan. 2021), p. 4948. DOI: 10.3390/app11114948. URL: <https://www.mdpi.com/2076-3417/11/11/4948> (visited on 07/10/2021).
- [Che+21] Dezhi Chen et al. “Mean Field Deep Reinforcement Learning for Fair and Efficient UAV Control”. In: *IEEE Internet of Things Journal* 8.2 (Jan. 2021), pp. 813–828. ISSN: 2327-4662. DOI: 10.1109/JIOT.2020.3008299.

- URL: <https://ieeexplore.ieee.org/abstract/document/9137257> (visited on 06/25/2021).
- [CS16] Davide Ceraso and Giandomenico Spezzano. “Controlling Swarms of Medical Nanorobots Using CPPSO on a GPU”. In: *2016 International Conference on High Performance Computing Simulation (HPCS)*. July 2016, pp. 58–65. DOI: 10.1109/HPCSim.2016.7568316. URL: <https://ieeexplore.ieee.org/document/7568316> (visited on 07/20/2021).
- [Dre21] Daniel S. Drew. “Multi-Agent Systems for Search and Rescue Applications”. In: *Current Robotics Reports* 2.2 (June 2021), pp. 189–200. ISSN: 2662-4087. DOI: 10.1007/s43154-021-00048-3. URL: <https://doi.org/10.1007/s43154-021-00048-3> (visited on 07/20/2021).
- [EH01] M. Egerstedt and Xiaoming Hu. “Formation Constrained Multi-Agent Control”. In: *IEEE Transactions on Robotics and Automation* 17.6 (Dec. 2001), pp. 947–951. ISSN: 2374-958X. DOI: 10.1109/70.976029. URL: <https://ieeexplore.ieee.org/document/976029> (visited on 07/05/2021).
- [Eng+19] Logan Engstrom et al. “Implementation Matters in Deep RL: A Case Study on PPO and TRPO”. In: Sept. 2019. URL: <https://openreview.net/forum?id=r1etN1rtPB> (visited on 07/28/2021).
- [Foe+16] Jakob Foerster et al. “Learning to Communicate with Deep Multi-Agent Reinforcement Learning”. In: *Advances in Neural Information Processing Systems* 29 (2016). URL: <https://proceedings.neurips.cc/paper/2016/hash/c7635bfd99248a2cdef8249ef7bfbe4-Abstract.html> (visited on 06/25/2021).
- [Foe+17] Jakob Foerster et al. “Counterfactual Multi-Agent Policy Gradients”. In: *arXiv:1705.08926 [cs]* (Dec. 2017). URL: <http://arxiv.org/abs/1705.08926> (visited on 07/18/2021).
- [GEK17] Jayesh K. Gupta, Maxim Egorov, and Mykel Kochenderfer. “Cooperative Multi-agent Control Using Deep Reinforcement Learning”. In: *Autonomous Agents and Multiagent Systems*. Ed. by Gita Sukthankar and Juan A. Rodriguez-Aguilar. Lecture Notes in Computer Science. Cham: Springer International Publishing, 2017, pp. 66–83. ISBN: 9783319716824. DOI: 10.1007/978-3-319-71682-4_5. URL: https://link.springer.com/chapter/10.1007/978-3-319-71682-4_5 (visited on 06/25/2021).
- [GHN19] Gregor H. W. Gebhardt, Maximilian Hüttenrauch, and G. Neumann. *Using M-Embeddings to Learn Control Strategies for Robot Swarms*. 2019. URL: <https://www.semanticscholar.org/paper/Using-M-Embeddings-to-Learn-Control-Strategies-for-Gebhardt-H%C3%BCttenrauch/9f550815f8858e7c4c8aef23665fa5817884f1b3> (visited on 06/20/2021).
- [Haa+18] Tuomas Haarnoja et al. “Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor”. In: *arXiv:1801.01290 [cs, stat]* (Aug. 2018). URL: <http://arxiv.org/abs/1801.01290> (visited on 07/20/2021).

- [HŠN19] Maximilian Hüttenrauch, Adrian Šošić, and Gerhard Neumann. “Deep Reinforcement Learning for Swarm Systems”. In: *Journal of Machine Learning Research* 20.54 (2019), pp. 1–31. ISSN: 1533-7928. URL: <http://jmlr.org/papers/v20/18-476.html> (visited on 06/14/2021).
- [Hu+20] Junyan Hu et al. “Occlusion-Based Coordination Protocol Design for Autonomous Robotic Shepherd Tasks”. In: *IEEE Transactions on Cognitive and Developmental Systems* (2020), pp. 1–1. ISSN: 2379-8939. DOI: 10.1109/TCDS.2020.3018549. URL: <https://ieeexplore.ieee.org/abstract/document/9173524> (visited on 06/20/2021).
- [IS19] Shariq Iqbal and Fei Sha. “Actor-Attention-Critic for Multi-Agent Reinforcement Learning”. In: *International Conference on Machine Learning*. PMLR, May 2019, pp. 2961–2970. URL: <http://proceedings.mlr.press/v97/iqbal19a.html> (visited on 07/09/2021).
- [KL51] S. Kullback and R. A. Leibler. “On Information and Sufficiency”. In: *The Annals of Mathematical Statistics* 22.1 (Mar. 1951), pp. 79–86. ISSN: 0003-4851. DOI: 10.1214/aoms/1177729694. URL: <http://projecteuclid.org/euclid.aoms/1177729694> (visited on 07/20/2021).
- [KMP13] Soumya Kar, José M. F. Moura, and H. Vincent Poor. “Q D-learning: A Collaborative Distributed Strategy for Multi-Agent Reinforcement Learning Through $\{\text{Rm Consensus}\} + \{\text{Rm Innovations}\}$ ”. In: *IEEE Transactions on Signal Processing* 61.7 (2013), pp. 1848–1862. DOI: 10.1109/TSP.2013.2241057.
- [LL07] Jean-Michel Lasry and Pierre-Louis Lions. “Mean Field Games”. In: *Japanese Journal of Mathematics* 2.1 (Mar. 2007), pp. 229–260. ISSN: 1861-3624. DOI: 10.1007/s11537-007-0657-8. URL: <https://doi.org/10.1007/s11537-007-0657-8> (visited on 07/10/2021).
- [Low+17] Ryan Lowe et al. “Multi-Agent Actor-Critic for Mixed Cooperative-Competitive Environments”. In: *Proceedings of the 31st International Conference on Neural Information Processing Systems*. NIPS’17. Long Beach, California, USA: Curran Associates Inc., Dec. 2017, pp. 6382–6393. ISBN: 9781510860964. URL: <https://dl.acm.org/doi/10.5555/3295222.3295385> (visited on 07/12/2021).
- [LT20] Xiangyu Liu and Ying Tan. “Attentive Relational State Representation in Decentralized Multiagent Reinforcement Learning”. In: *IEEE Transactions on Cybernetics* (2020), pp. 1–13. ISSN: 2168-2275. DOI: 10.1109/TCYB.2020.2979803. URL: <https://ieeexplore.ieee.org/document/9049415> (visited on 06/22/2021).
- [MA18] Igor Mordatch and Pieter Abbeel. “Emergence of Grounded Compositional Language in Multi-Agent Populations”. In: *arXiv:1703.04908 [cs]* (July 2018). URL: <http://arxiv.org/abs/1703.04908> (visited on 06/14/2021).

- [MLL07] Laetitia Matignon, Guillaume J. Laurent, and Nadine Le Fort-Piat. “Hysteretic Q-learning : An Algorithm for Decentralized Reinforcement Learning in Cooperative Multi-Agent Teams”. In: *2007 IEEE/RSJ International Conference on Intelligent Robots and Systems*. Oct. 2007, pp. 64–69. DOI: 10.1109/IRoS.2007.4399095. URL: <https://ieeexplore.ieee.org/document/4399095> (visited on 07/10/2021).
- [OA16] Frans A. Oliehoek and Christopher Amato. *A Concise Introduction to Decentralized POMDPs*. SpringerBriefs in Intelligent Systems. Springer International Publishing, 2016. ISBN: 9783319289274. DOI: 10.1007/978-3-319-28929-8. URL: <https://www.springer.com/gp/book/9783319289274> (visited on 07/21/2021).
- [Oli12] Frans A. Oliehoek. “Decentralized POMDPs”. In: *Reinforcement Learning: State-of-the-Art*. Ed. by Marco Wiering and Martijn van Otterlo. Adaptation, Learning, and Optimization. Berlin, Heidelberg: Springer, 2012, pp. 471–503. ISBN: 9783642276453. DOI: 10.1007/978-3-642-27645-3_15. URL: https://doi.org/10.1007/978-3-642-27645-3_15 (visited on 06/25/2021).
- [Ope+19] OpenAI et al. “Learning Dexterous In-Hand Manipulation”. In: *arXiv:1808.00177 [cs, stat]* (Jan. 2019). URL: <http://arxiv.org/abs/1808.00177> (visited on 07/20/2021).
- [Ott+20] Fabian Otto et al. “Differentiable Trust Region Layers for Deep Reinforcement Learning”. In: Sept. 2020. URL: <https://openreview.net/forum?id=qYZD-A01Vn> (visited on 06/14/2021).
- [Raf+19] Antonin Raffin et al. *Stable Baselines3*. 2019.
- [Raf+20] Antonin Raffin et al. *PPO — Stable Baselines3 1.1.0a11 Documentation*. 2020. URL: <https://stable-baselines3.readthedocs.io/en/master/modules/ppo.html#results> (visited on 06/14/2021).
- [Ras04] Carl Edward Rasmussen. “Gaussian Processes in Machine Learning”. In: *Advanced Lectures on Machine Learning: ML Summer Schools 2003, Canberra, Australia, February 2 - 14, 2003, Tübingen, Germany, August 4 - 16, 2003, Revised Lectures*. Ed. by Olivier Bousquet, Ulrike von Luxburg, and Gunnar Rätsch. Lecture Notes in Computer Science. Berlin, Heidelberg: Springer, 2004, pp. 63–71. ISBN: 9783540286509. DOI: 10.1007/978-3-540-28650-9_4. URL: https://doi.org/10.1007/978-3-540-28650-9_4 (visited on 07/10/2021).
- [San17] Andrew W. Sanders. *Drone Swarms*. May 2017. URL: <https://apps.dtic.mil/sti/citations/AD1039921> (visited on 07/20/2021).
- [Sch+15] John Schulman et al. “Trust Region Policy Optimization”. In: *International Conference on Machine Learning*. PMLR, June 2015, pp. 1889–1897. URL: <http://proceedings.mlr.press/v37/schulman15.html> (visited on 07/20/2021).

- [Sch+17] John Schulman et al. “Proximal Policy Optimization Algorithms”. In: *arXiv:1707.06347 [cs]* (Aug. 2017). URL: <http://arxiv.org/abs/1707.06347> (visited on 06/14/2021).
- [Sch+18] John Schulman et al. “High-Dimensional Continuous Control Using Generalized Advantage Estimation”. In: *arXiv:1506.02438 [cs]* (Oct. 2018). URL: <http://arxiv.org/abs/1506.02438> (visited on 07/21/2021).
- [Sil+17] David Silver et al. “Mastering the Game of Go Without Human Knowledge”. In: *Nature* 550.7676 (Oct. 2017), pp. 354–359. ISSN: 1476-4687. DOI: 10.1038/nature24270. URL: <https://www.nature.com/articles/nature24270> (visited on 06/25/2021).
- [Šoš+17] Adrian Šošić et al. “Inverse Reinforcement Learning in Swarm Systems”. In: *Proceedings of the 16th Conference on Autonomous Agents and Multi-Agent Systems. AAMAS ’17. São Paulo, Brazil: International Foundation for Autonomous Agents; Multiagent Systems, May 2017*, pp. 1413–1421. URL: <https://dl.acm.org/doi/10.5555/3091125.3091320> (visited on 07/21/2021).
- [TD20] Alexey Turchin and David Denkenberger. “Classification of Global Catastrophic Risks Connected with Artificial Intelligence”. In: *AI & SOCIETY* 35.1 (Mar. 2020), pp. 147–163. ISSN: 1435-5655. DOI: 10.1007/s00146-018-0845-5. URL: <https://doi.org/10.1007/s00146-018-0845-5> (visited on 07/20/2021).
- [Tea+21] Open-Ended Learning Team et al. “Open-Ended Learning Leads to Generally Capable Agents”. In: *arXiv:2107.12808 [cs]* (July 2021). URL: <http://arxiv.org/abs/2107.12808> (visited on 07/28/2021).
- [Ter+21] J. K. Terry et al. “PettingZoo: Gym for Multi-Agent Reinforcement Learning”. In: *arXiv:2009.14471 [cs, stat]* (June 2021). URL: <http://arxiv.org/abs/2009.14471> (visited on 07/09/2021).
- [Vas+17] Ashish Vaswani et al. “Attention Is All You Need”. In: *Advances in Neural Information Processing Systems* 30 (2017). URL: <https://papers.nips.cc/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html> (visited on 06/14/2021).
- [Vil09] Cédric Villani. *Optimal Transport: Old and New*. Grundlehren der mathematischen Wissenschaften. Berlin Heidelberg: Springer-Verlag, 2009. ISBN: 9783540710493. DOI: 10.1007/978-3-540-71050-9. URL: <https://www.springer.com/gp/book/9783540710493> (visited on 07/20/2021).
- [Vol+20] Michael Volpp et al. “Bayesian Context Aggregation for Neural Processes”. In: Sept. 2020. URL: <https://openreview.net/forum?id=ufZN2-aehFa> (visited on 06/14/2021).
- [Wai+19] Hoi-To Wai et al. “Multi-Agent Reinforcement Learning via Double Averaging Primal-Dual Optimization”. In: *arXiv:1806.00877 [cs, math, stat]* (Jan. 2019). URL: <http://arxiv.org/abs/1806.00877> (visited on 06/25/2021).

- [WKA17] Markus Waibel, Bill Keays, and Federico Augugliaro. *Drone Shows: Creative Potential and Best Practices*. 2017. DOI: 10.3929/ethz-a-010831954. URL: <https://www.research-collection.ethz.ch/handle/20.500.11850/125498> (visited on 07/20/2021).
- [Yan+18] Yaodong Yang et al. “Mean Field Multi-Agent Reinforcement Learning”. In: *International Conference on Machine Learning*. PMLR, July 2018, pp. 5571–5580. URL: <http://proceedings.mlr.press/v80/yang18d.html> (visited on 06/14/2021).
- [Zha+18] Kaiqing Zhang et al. “Fully Decentralized Multi-Agent Reinforcement Learning with Networked Agents”. In: *International Conference on Machine Learning*. PMLR, July 2018, pp. 5872–5881. URL: <http://proceedings.mlr.press/v80/zhang18n.html> (visited on 06/25/2021).
- [Zhe+18] Lianmin Zheng et al. “MAgent: A Many-Agent Reinforcement Learning Platform for Artificial Collective Intelligence”. In: *Proceedings of the AAAI Conference on Artificial Intelligence* 32.1 (Apr. 2018). ISSN: 2374-3468. URL: <https://ojs.aaai.org/index.php/AAAI/article/view/11371> (visited on 06/25/2021).
- [Zho+16] Zhengyuan Zhou et al. “Cooperative Pursuit with Voronoi Partitions”. In: *Automatica* 72 (Oct. 2016), pp. 64–72. ISSN: 0005-1098. DOI: 10.1016/j.automatica.2016.05.007. URL: <https://www.sciencedirect.com/science/article/pii/S0005109816301911> (visited on 07/05/2021).
- [ZYB21] Kaiqing Zhang, Zhuoran Yang, and Tamer Başar. “Multi-Agent Reinforcement Learning: A Selective Overview of Theories and Algorithms”. In: *arXiv:1911.10635 [cs, stat]* (Apr. 2021). URL: <http://arxiv.org/abs/1911.10635> (visited on 06/25/2021).

Appendix A.

Experiment Hyper-Parameters and Overview

Experiment	Entities	Dynamics	Batch size	Environment steps per training step	Max Training Steps
Rendezvous	20 agents	Double-Integrator Unicycle	1000	164000	160
Single-Evader Pursuit	10 agents, 1 evader	Single-Integrator Unicycle	10200	102000	500
Multi-Evader Pursuit	50 agents, 5 evaders (respawning)	Single-Integrator Unicycle	10200	102000	500
Box Assembly	10 agents, 4 boxes	Single-Integrator Unicycle	5000	250000	200
Box Clustering (2 clusters)	10 agents, 2 clusters with 2 boxes each	Single-Integrator Unicycle	10000	512000	2000
Box Clustering (3 clusters)	20 agents, 3 clusters with 4 boxes each	Single-Integrator Unicycle	100000	12800000	500