

A Model-Based Solution to the Offline Multi-Agent Reinforcement Learning Coordination Problem

Paul Barde
Mila - Quebec AI Institute
McGill University
Montreal, Canada
bardepau@mila.quebec

Jakob Foerster
University of Oxford
Oxford, United Kingdom

Derek Nowrouzezahrai
Mila - Quebec AI Institute
McGill University
Montreal, Canada

Amy Zhang
The University of Texas at Austin
Meta AI - FAIR
Austin, USA

ABSTRACT

Training multiple agents to coordinate is an essential problem with applications in robotics, game theory, economics, and social sciences. However, most existing Multi-Agent Reinforcement Learning (MARL) methods are online and thus impractical for real-world applications in which collecting new interactions is costly or dangerous. While these algorithms should leverage offline data when available, doing so gives rise to what we call *the offline coordination problem*. Specifically, we identify and formalize the *strategy agreement (SA)* and the *strategy fine-tuning (SFT)* coordination challenges, two issues at which current offline MARL algorithms fail. Concretely, we reveal that the prevalent model-free methods are severely deficient and cannot handle coordination-intensive offline multi-agent tasks in either toy or MuJoCo domains. To address this setback, we emphasize the importance of inter-agent interactions and propose the very first model-based offline MARL method. Our resulting algorithm, Model-based Offline Multi-Agent Proximal Policy Optimization (MOMA-PPO) generates synthetic interaction data and enables agents to converge on a strategy while fine-tuning their policies accordingly. This simple model-based solution solves the coordination-intensive offline tasks, significantly outperforming the prevalent model-free methods even under severe partial observability and with learned world models.

KEYWORDS

Multi-Agent Learning; Coordination; Offline Reinforcement Learning; Model-Based Reinforcement Learning; Deep Learning; World Model

ACM Reference Format:

Paul Barde, Jakob Foerster, Derek Nowrouzezahrai, and Amy Zhang. 2024. A Model-Based Solution to the Offline Multi-Agent Reinforcement Learning Coordination Problem. In *Proc. of the 23rd International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2024)*, Auckland, New Zealand, May 6 – 10, 2024, IFAAMAS, 17 pages.



This work is licensed under a Creative Commons Attribution International 4.0 License.

Proc. of the 23rd International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2024), N. Alechina, V. Dignum, M. Dastani, J.S. Sichman (eds.), May 6 – 10, 2024, Auckland, New Zealand. © 2024 International Foundation for Autonomous Agents and Multiagent Systems (www.ifaamas.org).

1 INTRODUCTION

Multi-agent problems are ubiquitous in real-world scenarios, including traffic control, distributed energy management, multi-robot coordination, auctions, marketplaces, and social networks [6, 8, 10, 13, 15, 16, 21, 25, 42, 44, 53, 58]. This makes the development of efficient multi-agent algorithms a crucial research area in artificial intelligence and machine learning with substantial implications in various fields including robotics, game theory, economics, and social sciences [35, 37, 51, 60, 68]. However, existing methods are mostly online and require interacting with the environment throughout learning which often makes them costly or even dangerous for real-world applications [34]. In contrast, offline Reinforcement Learning (offline RL) obviates the need for interactions with the environment as it allows learning from existing datasets that do not have to be collected by experts. It is therefore well suited to tasks where: a) we cannot afford to materialize the situation in practice, b) it is unfeasible to build a simulator, and c) there exist datasets of realizations of such situations. Consequently, we hypothesize that offline multi-agent approaches will be key for tackling real-world multi-agent problems. Let us imagine for instance trying to understand how autonomous actors (i.e., governments, international organizations, industries, etc.) must maneuver to reduce the severity of a worldwide pandemic while preventing economic collapse. It goes without saying that: a) starting pandemics is not a viable way to gain real-world practice; b) building a simulator is a colossal task that would suggest emulating our society and its economy; and c) the impact of past decisions (such as implementing lockdown policies, travel restrictions, and vaccination campaigns) on the unfolding of the pandemics and the economy is well-documented. Hopefully, these records can be used to derive new strategies in the future.

The offline coordination challenge. In general, actors such as individuals, organizations, robots, software processes, or cars are self-governed and ultimately act autonomously. However, during the offline learning phase, it is reasonable to assume that learners can share pieces of information, which makes the Centralized Training Decentralized Execution (CTDE) formulation a natural approach for the offline training regime. Unfortunately, as shown in this work, simply combining CTDE MARL and offline RL methods on a dataset of multi-agent interactions does not necessarily yield policies that perform well together. This is what we call *the offline*

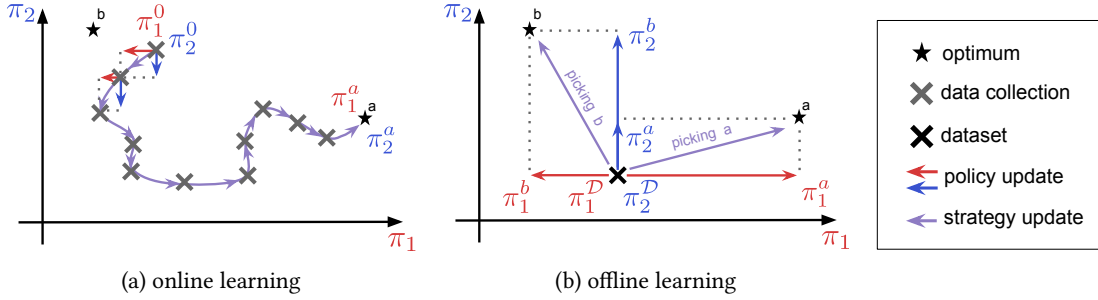


Figure 1: Comparing online learning and offline learning in policy space to illustrate the offline coordination problem. (a) During online learning, agents continuously interact using their current policies and collect new data that informs the next update of the co-evolution from (π_1^0, π_2^0) to (π_1^a, π_2^a) . (b) During offline learning, agents cannot collect new data and thus, they can only estimate updates from the dataset of interactions (here collected with π_1^D and π_2^D). To reach an optimal strategy agents must (1) agree on which optimum to target between \star^a or \star^b , – i.e. solve *strategy agreement* –, and (2) respectively derive the policy corresponding to that strategy $(\pi_1^{\star j}$ and $\pi_2^{\star j})$, – i.e. solve *strategy fine-tuning*.

coordination problem: 1) multi-agent solutions require agents to coordinate, that is to act coherently as a group such that individual behaviors combine into an efficient team strategy; 2) learning to coordinate is particularly challenging in the absence of interactions during training.

In this work, we propose that multi-agent coordination can be decomposed into two distinct challenges. First, since most multi-agent tasks allow multiple optimal team strategies [5] – some of which may only vary on how they break symmetries present in the task [24] – agents must collectively select one strategy over another such that they individually converge toward coherent behaviors. We refer to this coordination challenge as *Strategy Agreement (SA)*. Additionally, for a chosen team strategy, agents have to precisely calibrate and adjust their behaviors to one another in order to reach the corresponding optimal group behavior. We call this *Strategy Fine-Tuning (SFT)*. During online learning, agents continuously interact together in the environment, therefore, changes due to one agent’s local optimization directly impact other agents’ experiences and teammates are able to adapt: coordination occurs through interactive trial and error. Conversely, when learning from a fixed dataset of interactions, agents cannot interact with other agents to generate new data that explicitly measures the outcomes of current policies in the environment (i.e., which global strategy is being chosen and how individual behaviors blend together). Therefore, it is difficult for offline learners to coordinate. Figure 1 illustrates this in policy space: (a) During online learning, agents continuously interact to collect up-to-date data that inform the optimization of both individual policies and team strategy. This allows agents to converge to a global optimum by co-adapting and improving on each other’s updates. (b) In the offline setting however, agents must independently decide towards which of the two optima they aim to converge (i.e., strategy agreement between $\pi_i^{\star a}$ or $\pi_i^{\star b}$, $i = 1, 2$). Assuming they pick \star^a , they must then derive their corresponding optimal policy (i.e., strategy fine-tuning toward $\pi_i^{\star a}$, $i = 1, 2$). Crucially, offline agents must rely only on interactions present in the dataset (corresponding to policies π_i^D , $i = 1, 2$), thus likely without interactions corresponding to their past or current policies.

Current methods [27, 43, 64] deal with offline MARL by simply extending single-agent offline RL to the multi-agent setting. To do so, they either consider that agents are independent learners, or they leverage the CTDE paradigm. Provided that the dataset has enough coverage, it is in theory possible for CTDE agents to learn the different optimal value functions and policies. Yet, as we will see in Section 6, agents may still fail to agree on which team strategy to pick. We highlight such offline coordination failure in an offline version of the well-established Iterated Coordination Game [5]. The crux of *offline* MARL is therefore that learners cannot interact together in the environment to collect new data that measures how individual policies blend together and what are the outcomes of the current team strategy. This relates to the absence of interactions during learning and is not a centralization issue. Thus, we motivate the need for generating synthetic data that measure the outcomes of the team’s current strategy to allow for coordinated agents.

We propose MOMA-PPO, a simple model-based approach that generates synthetic interactions. We show that it allows for offline coordination even in complex Multi-Agent MuJoCo (MAMuJoCo) [45] tasks with partial observability and learned world model. Across tasks and domains, MOMA-PPO significantly outperforms the offline MARL baselines. Our contributions are:

- identifying and defining the offline coordination problem, an issue that has been overlooked by the offline MARL community;
- proposing benchmarks in the form of new datasets and partially observable multi-agent tasks that test for offline coordination;
- showing that current methods, which are all model-free, fail at offline coordination even in simple environments;
- suggesting a link between the offline coordination problem and the lack of inter-agent interactions during learning that is inherent to model-free offline approaches;
- proposing to address the coordination problem with the first-ever model-based offline MARL method;
- experimentally validating the benefits of coordinating by interacting through a world model.

2 BACKGROUND

We consider the formalism of Dec-POMDP [41] with states $s_t \in \mathcal{S}$. There are $i = 1, \dots, N_A$ agents that partially observe this state from their stochastic observation function $s_t^i \sim O^i(s^i|s_t)$ and choose action $a_t^i \in \mathcal{A}^i$ at every time step. Each agent has an action-observation history given by $h_t^i = \{s_0^i, a_0^i, \dots, s_t^i\}$ and picks its action a_t^i using its stochastic policy $\pi_i(a^i|h_t^i)$. The environment’s transition depends on the joint action a_t such that $s_{t+1} \sim P(s'|s_t, a_t)$. The game is fully cooperative so all agents receive the same reward and $r_t^i = r_t = r(s_t, a_t) \in \mathbb{R} \forall i$. The agent’s objective is therefore to maximize the expected team return $J = \mathbb{E}_\tau R(\tau)$ where $R(\tau) = \sum_t \gamma^t r_t$ with discount factor γ and $\tau = \{s_0, a_0, r_0, \dots, s_F\}$ is a trajectory with absorbing state s_F . Note that some trajectories can become arbitrarily long in which case we truncate them and use the value of the last state as an estimation of the return-to-go.

We adopt the CTDE framework [17, 38], where at training time, individual observations, policies, and value functions are available to all the agents. For simplicity, we assume access to the global state s_t and the observation functions O^i . Yet, this is not a requirement and it is always possible to define the global state as the concatenation of the agents’ observations, i.e., $s_t = \{s_t^1, \dots, s_t^{N_A}\}$. In most problems, which are special cases of Dec-POMDPs referred to as Dec-MDPs, such concatenation will fully observe the environment. At execution, agents must act independently and the joint policy is approximated by individual decentralized policies as $a_t \sim \pi(a|s_t) \approx \prod_{i=0}^{N_A} \pi_i(a^i|h_t^i)$.

Finally, we consider Offline Learning [20, 34], where agents have only access to a fixed dataset of trajectories \mathcal{D} and cannot collect additional interactions with the environment during learning.

3 METHOD

In this work, we propose MOMA-PPO, a Dyna-like [55] model-based approach to multi-agent CTDE offline learning that relies on PPO [50]. The method can be decomposed into two steps: 1) learning a world model from the dataset, and 2) using the world model to train the agents’ policies.

3.1 Learning a centralized world model ensemble

MOMA-PPO leverages the CTDE assumptions to learn centralized models that predict the next state, reward, and termination condition from the current state and actions. When learning in an approximate world model, RL agents might learn to exploit the world model’s reconstruction inaccuracies to reap more rewards in simulation, eventually producing incoherent behaviors that perform poorly in the real world [22]. One way to avoid this is to penalize the agents for going into regions of the state-action space where the world model is uncertain about its predictions [67]. Learning an ensemble of models enables estimating the world model’s epistemic uncertainty due to the limited amount of learning data in the offline dataset. Each model comprises two diagonal Gaussians $\mathcal{N}(\mu_T, \sigma_T^2)$ and $\mathcal{N}(\mu_r, \sigma_r^2)$ that respectively model the next state s' and the reward r . Models also predict whether or not the next state is terminal using a Bernoulli distribution $\text{Bern}(p_d)$. Distributions’ $\mu_T, \mu_r, \sigma_T, \sigma_r$, and p_d are parametrized by neural networks conditioned on

the current global state s and the joint action a . The parameters are learned from the offline dataset \mathcal{D} using Gaussian negative log-likelihood for $\mathcal{N}(\mu_T, \sigma_T^2)$ and $\mathcal{N}(\mu_r, \sigma_r^2)$, and binary cross-entropy for $\text{Bern}(p_d)$. In practice, we train $N_m = 7$ models and keep the best $N = 5$ based on their average validation accuracy across the next states and rewards [67]. We estimate the epistemic uncertainty of the reward using the variance of the predicted rewards across the ensemble:

$$\epsilon_r = \frac{\sum_{m=1}^N (\hat{r}_m - \bar{r})^2}{N-1}, \quad \bar{r} = \frac{\sum_{m=1}^N \hat{r}_m}{N}.$$

We also estimate the epistemic uncertainty of the general prediction by concatenating the next state and the reward and computing the Frobenius norm of the ensemble covariance matrix:

$$\epsilon_g = \|\text{cov}(x_i, x_j)\|_F, \quad \text{cov}(x_i, x_j) = \frac{\sum_{m=1}^N (\hat{x}_{i,m} - \bar{x}_i)(\hat{x}_{j,m} - \bar{x}_j)}{N-1},$$

where x_i and x_j are components of the vector resulting from the concatenation of the predicted next state vector \hat{s}' and the predicted reward scalar \hat{r} . At this point, we define a world model based on the ensemble such that:

$$\hat{s}_{t+1}, \bar{r}_t, \bar{f}_t, \epsilon_{r,t}, \epsilon_{g,t} \sim \mathcal{M}(\cdot|s_t, a_t),$$

where \bar{f}_t is a mask equal to 0, if the model predicts that we reached an absorbing state, and 1 otherwise. \bar{r}_t is the mean predicted reward across the ensemble and \bar{f}_t results from a majority vote between the members of the ensemble. The mean state would likely be out-of-distribution and lack the structure of real states which would impede learning and evaluation. Therefore, \hat{s}_{t+1} is instead sampled uniformly amongst each ensemble member’s next state. Finally, to avoid unrealistic values, \bar{r}_t and \hat{s}_{t+1} are clipped to the minimum bounding box of the offline dataset while uncertainties estimations are limited to a specified threshold l_ϵ :

$$\begin{aligned} \min_{r \in \mathcal{D}} r \leq \bar{r}_t \leq \max_{r \in \mathcal{D}} r \\ \min_{s_i \in \mathcal{D}} s_i \leq \hat{s}_{t+1,i} \leq \max_{s_i \in \mathcal{D}} s_i \\ \text{and } \epsilon_r \leq l_\epsilon, \epsilon_g \leq l_\epsilon. \end{aligned} \quad \forall i \in [0, q] \mid \mathcal{S} \subset \mathbb{R}^q,$$

3.2 Model-based Offline Multi-Agent PPO (MOMA-PPO)

Once \mathcal{M} has been trained on the offline dataset \mathcal{D} , it can be used to train online reinforcement learning algorithms in a Dyna-like manner. Here, we use MAPPO, a CTDE multi-agent version of PPO [65].

The synthetic data used to train the PPO policies is collected by sampling states from the offline dataset \mathcal{D} and using the current policies π_i alongside the world model \mathcal{M} to generate PPO’s training rollouts of size k . Terminating a rollout when the world model uncertainty exceeds l_ϵ allows for adaptive rollout length and avoids training the

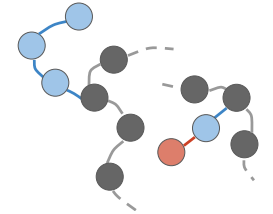


Figure 2: Model-based rollouts generation (blue) from dataset’s states (grey). Red denotes early termination and $k = 3$.

policies on unfeasible data. Rollouts are always terminated with a timeout mask ζ_t such that:

$$\zeta_t = 1 - \mathbb{I}(t = k - 1 \cup \epsilon_{g,t} \geq l_\epsilon),$$

where \mathbb{I} is the indicator function. Note that the timeout mask ζ_t is different from the model-predicted mask f_t which is solely related to the environment and indicates terminal absorbing states that occur for instance when the agent dies or reaches the goal. On the PPO side, we adapt the generalized advantage estimation [49] to account for the timeouts ζ_t and ensure that there is no accumulation across rollouts while computing returns. We use the value of the last state as an estimation of the return-to-go (see Appendix A).

The generation of length $k = 3$ rollouts is illustrated in Figure 2 where it can be seen that a rollout is interrupted early because the world model was unreliable when generating the associated state (shown in red).

In addition to the adaptive rollout length, the model epistemic uncertainty is used to penalize the agents so they avoid exploiting the world model in poorly reconstructed regions of the state space. The final uncertainty-penalized reward is given by:

$$\tilde{r}_t = \bar{r}_t - \lambda_r \epsilon_r - \lambda_g \epsilon_g,$$

where ϵ_r and ϵ_g are hyperparameters that weigh the severity of the penalty.

Algorithm 1 MOMA-PPO

Require: offline dataset \mathcal{D} , world model \mathcal{M} , rollout horizon k , rollout batch size b , uncertainty penalty coefficients λ_r and λ_g , uncertainty threshold l_ϵ , MAPPO agents.

Initialize MAPPO policies π_i and value function V .

for epoch $1, 2, \dots$ **do**

▷ Generate synthetic data

Initialize an empty rollout buffer $\mathcal{R} \leftarrow \emptyset$.

for $1, 2, \dots, b$ (in parallel) **do**

Sample history $h_t = \{h_t^i\}_{i=1}^{N_A}$ from \mathcal{D} .

for $j = t, t + 1, \dots, t + k - 1$ **do**

Sample $a_j^i \sim \pi_i(a^i | h_j^i) \forall i$.

Sample $\hat{s}_{j+1}, \hat{r}_j, \hat{f}_j, \epsilon_{r,j}, \epsilon_{g,j} \sim \mathcal{M}(\cdot | s_j, a_j)$.

Compute $\tilde{r}_j = \hat{r}_j - \lambda_r \epsilon_r - \lambda_g \epsilon_g$.

Compute $\zeta_j = 1 - \mathbb{I}(j = t + k - 1 \cup \epsilon_{g,j} \geq l_\epsilon)$

Add sample $(h_j, a_j, \tilde{r}_j, \hat{f}_j, \zeta_j, \hat{s}_{j+1})$ to \mathcal{R} .

Get h_{j+1}^i from h_j^i, a_j^i and \hat{s}_{j+1} . Set $s_{j+1} = \hat{s}_{j+1}$.

▷ Train agents

Use synthetic rollouts in \mathcal{R} to train policies π_i and value function V with MAPPO.

return multi-agent policies π_i .

Practical considerations and limitations. As stated in Section 2, we assume CTDE and that the global state s_t fully observes the environment, therefore we do not equip the world model with memory. The agents, on the other hand, only have access to partial observations and must rely on action-observation histories h_t^i . In practice, we restrict action-observation histories to 10 steps in the past (either from the dataset or from the generated rollouts) and process them with one layer of self-attention [59] followed by one layer of soft-attention [1]. The resulting embeddings are concatenated to the agent’s current state, and for simplicity, we abuse notation by denoting this “memory enhanced” state with h_t^i .

In MAPPO our centralized value function uses the QMIX value-decomposition [47] with $w^i(s_t)$ and $b(s_t)$ provided by a learnable neural network: $V(s_t) \triangleq \sum_i w^i(s_t) V^i(h_t^i) + b(s_t)$.

Finally, it is important to note that the task of the MOMA-PPO agents is quite different from the task of the agents that generated the dataset. First, the MOMA-PPO agents’ initial state distribution is now the dataset’s state distribution. Then the reward of the task has been altered to account for the model uncertainty. Last but not least, agents are only allowed to stray k steps away from the dataset’s coverage. While this restriction mitigates world model abuse, it can also prevent the agents from discovering goals that are further away from the offline data. Our resulting model-based offline multi-agent method is illustrated in Algorithm 1 and more details are provided in Appendix A.

4 RELATED WORK

Multi-agent coordination. Coordination has been a challenge of interest since the early works on cooperative MARL [5, 7, 11, 12, 36] and has consistently been a central focus of the multi-agent literature [26, 33, 38, 69]. While different works consider different aspects of coordination – such as behavior predictability and synchronous sub-policy selection [48], structured team exploration [40] or the emergence of communication and cooperative guiding [2, 32, 63] – our definition of coordination is closest to the seminal work of [5]. Indeed, we consider coordination in terms of agents agreeing to individually follow the same team strategy (that is a policy over joint actions) and finetuning their behavior to one another in tasks where multiple distinct optimal team strategies exist. A similar notion of coordination has been used in the *zero-shot coordination problem* investigated by [24] where agents are trained so that they are able to perform with agents they have never seen before. Yet, while their focus is on deriving standardized coordinated strategies that can generalize to unseen teammates, coordination is still learned through online interactions.

Coordination with teammates without direct interactions is often referred to as *ad-hoc coordination* [4, 54]. Recent works assume access to offline demonstrations of the teammates’ behaviors. These can be used to guide the agent’s self-play training toward adopting the appropriate equilibrium (or “social conventions”) of its future teammates [33, 57]. Similarly, [9] uses offline data to learn a model of the teammate’s behavior and use it to train the agent to coordinate with that ally. In ad-hoc coordination, teammates’ behaviors are fixed and can be estimated a priori from the dataset: the learner merely has to identify the team strategy and adopt it. Conversely, in our offline coordination setting, all the agents are learning and therefore have unknown changing behaviors: they must identify the different potential team strategies and agree on which one to follow (*strategy agreement*). Simultaneously, they must finetune their behaviors to one another in order to reach this team policy (*strategy fine-tuning*). All this without being able to interact with other agents or the environment.

Offline MARL. Recent works have been investigating offline solutions to the MARL problem. All of these methods build on model-free single-agent approaches and constrain the policy to stay in the dataset’s distribution by using either SARSA-like schemes

(such as ICQ [64] and IQL [30]) or policy regularization (such as CQL [31] and TD3+BC [19]).

Some methods investigate specific modifications to improve performance in the multi-agent setting. For instance, in the decentralized setting, MABCQ [27] enforces an optimistic importance sampling modification that assumes that independent agents will strive toward similar high-rewarding states, yet since this does not discriminate between which high-rewarding state to favor, the strategy agreement issue remains. For discrete actions spaces problems, [56] propose a Transformer-based approach that learns a centralized teacher and distills its policy into independent student policies. Finally, OMAR [43] proposes to alleviate miscoordination failure in offline MARL by adding a zeroth order optimization method on top of multi-agent CQL, achieving state-of-the-art performance on a variety of tasks. We share these works’ goal of learning coordinated and efficient multi-agent teams in the offline setting. Yet, we believe that interacting learners and agents are essential to coordination and therefore take a different approach by focusing on model-based methods rather than model-free ones.

Offline model-based RL. Model-based approaches have been investigated in the single-agent offline RL setting. Notoriously, MOPO [67] proposed to learn an ensemble-based world model and use it to generate rollouts from the offline dataset to train a SAC agent [23]. They also proposed an uncertainty-based reward penalty to prevent the learner from exploiting the model. MOREL [28] takes a similar approach but prevents model abuse by learning a pessimistic MDP in which states that are outside of the dataset coverage become absorbing terminal states. COMBO [66] proposed a similar but more conservative version of MOPO by using CQL instead of SAC and learning on both generated and dataset’s states. Finally, ROMI [61] also uses model-free offline RL to derive a policy from a model-based augmented offline dataset, yet they enforce additional conservatism by learning a reverse policy and dynamics model to generate rollouts that lead to target states contained in the dataset. This mitigates against generating rollouts outside of the dataset’s coverage.

We believe that offline RL algorithms are ill-suited to learn on non-stationary data such as the one generated by updating policies be it in a world model or in a real environment. Therefore, our method MOMA-PPO uses an online RL algorithm instead of an offline one. Additionally, to enforce conservatism and avoid world model exploitation we use both uncertainty penalty and early rollout termination. In that sense, MOMA-PPO unifies what would be multi-agent extensions of MOPO and MOREL. Yet, instead of penalizing aleatoric uncertainty as in MOPO, we focus on epistemic uncertainty and estimate it by monitoring the coherence between the different models in the ensemble. Section 6.3 reports the benefits of using the epistemic uncertainty (due to the finite amount of training data) over the aleatoric uncertainty (from the environment’s stochasticity). Also, MOMA-PPO’s early rollout termination is done with timeouts rather than with MOREL’s terminal states meaning that it does not penalize agents while still avoiding using unfeasible rollouts to train them. Finally, in contrast with MOREL and MOPO which use off-policy algorithms, MOMA-PPO is based on MAPPO, an on-policy method that has reliably achieved state-of-the-art performance in MARL tasks [65]. It is therefore well suited to the

multi-agent setting in which slight changes in a teammate’s policy can drastically impact the overall group behavior and quickly make previous interactions obsolete.

Model-based multi-agent RL. In the online setting, model-based approaches aim to improve sample efficiency by reducing the number of interactions with the environment. Therefore, these methods use off-policy schemes and focus on how and when to collect additional environment data for refining the world model and the policies [62, 70, 71]. In offline MARL, sample efficiency is not a consideration since the environment data has already been collected offline and additional samples are not an option. Yet, as we show here, model-based approaches can benefit multi-agent coordination in the offline setting by allowing multiple learners to interact through the world model.

5 BASELINES, ENVIRONMENTS, TASKS, AND DATASETS

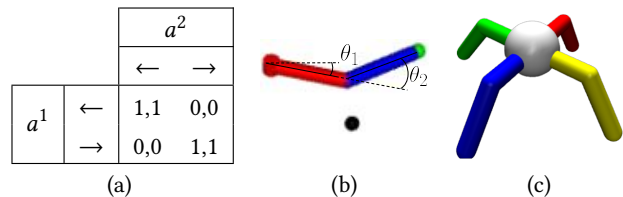


Figure 3: Environments illustrations. (a) Pay off matrix of the Iterated Coordination Game. (b) Two-agent Reacher, red and blue agents respectively control the torque on θ_1 and θ_2 . (c) Four-agent Ant, each agent controls a different limb (shown with different colors). In PO tasks, agents only observe the limb they control while the torso observations – in white – are available only to the yellow agent.

5.1 Baselines

We compare with a large and varied array of baselines. First, we consider a centralized training *and centralized execution* approach by observing the global state *and selecting the joint action*. We use IQL [30], a state-of-the-art single-agent model-free offline RL algorithm for this setting. Considering centralized execution gives an upper bound on what can be achieved in terms of strategy agreement. Indeed, a single learner controls all the agents and can thus choose for the whole team the strategy to adopt.

Then, we extend IQL [30] to the multi-agent setting by using the QMIX value decomposition on the Q and V value networks. This gives MAIQL, a very competitive model-free CTDE offline MARL algorithm that should allow for fine-tuning with additional online data after training. We refer to the finetuned version as MAIQL-ft and follow IQL’s finetuning procedure [30]: MAIQL-ft is first trained to convergence on the offline data and then finetuned for the same number of training steps by progressively introducing additional interaction data. At the end of finetuning, the replay buffer contains as many offline interactions as online ones. For completeness, we also consider MATD3+BC the CTDE multi-agent version of [19] with QMIX.

Recent literature in model-free MARL [14, 39] and notably in the offline setting [43], advocates for decentralized value functions and independent learners. Therefore, we consider several independent learner approaches, starting with Independent Behavioral Cloning (IBC). Despite its simplicity, BC [46] produces surprisingly efficient baselines for Imitation Learning and offline RL [3, 52]. Finally, we consider the independent learners extensions to [31] and [19], respectively ICQL and ITD3+BC, as well as the state-of-the-art model-free offline MARL method, OMAR [43].

5.2 Offline Iterated Coordination Game

To illustrate the strategy agreement coordination challenge, we propose an offline version of the Iterated Coordination Game [5] presented in Figure 3 (a). Agents must pick the same direction in order to succeed and we investigate three offline datasets of interactions (see Table 1). In the most favorable dataset, data is collected by coordinated agents that select the same option of going right most of the time. In the less favorable setting, agent 1 goes left most of the time while agent 2 is more likely to go right. In the neutral setting, agents act uniformly. However, each dataset does contain both coordinated and uncoordinated behaviors in which agents simultaneously choose the same – respectively, different – directions.

It is therefore straightforward for a centralized critic to learn that $Q(\rightarrow, \rightarrow) = Q(\leftarrow, \leftarrow) = 1$ while $Q(\rightarrow, \leftarrow) = Q(\leftarrow, \rightarrow) = 0$, regardless of the specific dataset. Yet, decentralized actors remain unaware of whether they should go left or right since both strategies are equivalent and actors have no way of consistently picking one over the other.

Table 1: Policies used to collect the datasets in the Iterated Coordination Game and resulting average scores of the datasets

	$P(a^1 = \rightarrow)$	$P(a^2 = \rightarrow)$	Avg. Score
favorable	0.75	0.75	0.623
neutral	0.5	0.5	0.502
unfavorable	0.25	0.75	0.375

5.3 Offline MAMuJoCo

Building upon D4RL [18] and MAMuJoCo [45], we propose offline multi-agent continuous control tasks with various datasets and full or partial observability.

Two-agent Reacher with a mixture-of-expert dataset. To investigate strategy agreement in a more complex continuous control setting, we propose a two-agent version of the Reacher environment as shown in Figure 3 (b). The offline dataset is collected as follows: in the first stage, we train online MAPPO on the fully observable two-agent Reacher task (every agent observes all the joint angles and velocities as well as the target position – in black – and the target to fingertip – in green – vector). Depending on the seed of the run, teams converge to counter-clockwise ($\theta_2 \geq 0$ as in Figure 3 (b)) or clockwise ($\theta_2 \leq 0$) arm bends. Thus we can build a mixture-of-expert dataset by combining equal proportions of demonstrations from clockwise and counter-clockwise teams. Finally, we explore the impact of Full Observability (FO) versus Partial Observability (PO) by considering three types of observation

functions: *all-observant* (FO: every agent fully observes the environment), *independent* (PO: each agent only sees the target and the velocity and angle of the joint it controls), and *leader-only* (PO: both agents observe the two joints but only the red agent observes the target’s position). Note that with PO no agent observes the target to fingertip vector. Thus, in *leader-only*, only the leader can estimate whether or not the fingertip matches the target. In *independent*, no agent has enough information to estimate this, yet both observe the target position. These tasks are very challenging and require agents to agree on following a specific convention (either clockwise or counter-clockwise arm bend) to reach a given target location and get rewards.

Four-agent Ant. Similarly, we use a MAMuJoCo-like decomposition [45] of the D4RL [18] offline ant task to make it multi-agent: each individual limb (composed of two joints) is controlled by a different agent. For the offline datasets, we use the single-agent D4RL datasets and consider two types of observation functions. For fully observable tasks, every agent observes the whole robot: torso observations (i.e., vertical position, orientation, angular and translational velocities) and the observations of all the limbs (i.e., angle and angular velocity of each joint). For the partially observable tasks, agents only observe the limb they control, and the torso observations are only made available to the yellow limb agent. PO tasks are very challenging in this case because only the yellow agent knows if the ant is moving in the correct direction and it must therefore learn to “steer” the whole robot.

Table 2 details all the datasets’ score distributions as well as the reference agents’ returns – expert and random – that are used to normalize performance. Note that ant datasets are from D4RL [18] (single agent datasets that we split into multi-agent observations and actions) while we generated the two-agent reacher mixture-of-expert dataset using MAPPO.

Table 2: Normalized measures of datasets’ scores distributions and normalization performances.

Scores		min	mean	median	max	expert	random
		(%)	(%)	(%)	(%)	return	return
reacher	expert-mix	38.8	100.0	98.9	152.3	-4.237	-11.145
ant	random	-3	6.4	7.2	10.3	3879.7	-325.6
	medium	-4.8	80.2	95.1	107.2		
	full-replay	-22.4	72.0	77.8	134.3		
	expert	-32.8	117.4	129.4	142.5		

6 RESULTS

Table 3: Teams’ performances on the Iterated Coordination Game. MOMA-PPO is the only decentralized execution method to solve it for all the datasets.

	IQL	MAIQL	IBC	MOMA-PPO
fav.	1. ± 0.	1. ± 0.	1. ± 0.	1. ± 0.
neutral	1. ± 0.	0.9 ± 0.1	0.55 ± 0.11	1. ± 0.
unfav.	1. ± 0.	0. ± 0.	0. ± 0.	1. ± 0.

Table 4: Teams’ performances on two-agent Reacher with mixture-of-experts dataset for different observation functions. Scores are normalized with expert and random performances. (a) Independent learners fail on datasets that contain a mixture of incompatible experts while MOMA-PPO (and to some extent MAIQL) are able to coordinate agents. (b) Current model-free methods are unable to adapt agents’ behaviors while MOMA-PPO significantly outperforms the baseline across all settings.

(a) Two-agent Reacher, mixture-of-experts dataset and full/partial observability

Algorithms		model-free							model-based (ours)
		centralized	CTDE		independent learners				CTDE
		IQL	MAIQL	MATD3+BC	IBC	ITD3+BC	ICQL	IOMAR	MOMA-PPO
FO	all-observant	1.07 ± 0.01	0.96 ± 0.05	1.04 ± 0.01	1.02 ± 0.01	0.78 ± 0.00	0.48 ± 0.06	0.73 ± 0.01	1.07 ± 0.01
PO	independent		0.92 ± 0.04	0.59 ± 0.03	0.76 ± 0.04	0.30 ± 0.11	0.46 ± 0.04	0.45 ± 0.02	0.95 ± 0.06
	leader-only		0.80 ± 0.05	0.73 ± 0.02	0.84 ± 0.02	0.48 ± 0.04	0.31 ± 0.05	0.39 ± 0.02	1.00 ± 0.01

(b) Four-agent Ant, various datasets and full/partial observability

Algorithms		model-free							model-based (ours)
		centralized	CTDE		independent learners				CTDE
		IQL	MAIQL	MAIQL-ft	IBC	ITD3+BC	ICQL	IOMAR	MOMA-PPO
FO	ant-random	0.12 ± 0.00	0.28 ± 0.01	0.28 ± 0.03	0.31 ± 0.00	0.22 ± 0.02	0.08 ± 0.00	0.08 ± 0.00	0.52 ± 0.07
	ant-medium	0.97 ± 0.02	0.85 ± 0.02	0.81 ± 0.02	0.84 ± 0.01	1.04 ± 0.00	0.88 ± 0.12	1.10 ± 0.03	1.29 ± 0.06
	ant-full-replay	1.22 ± 0.02	0.77 ± 0.21	0.95 ± 0.13	1.20 ± 0.01	1.33 ± 0.01	1.21 ± 0.02	1.30 ± 0.00	1.42 ± 0.07
	ant-expert	1.26 ± 0.01	1.24 ± 0.00	1.06 ± 0.07	1.24 ± 0.00	1.25 ± 0.02	0.73 ± 0.15	1.16 ± 0.01	1.49 ± 0.01
PO	ant-random		0.31 ± 0.00	0.34 ± 0.04	0.31 ± 0.00	0.31 ± 0.00	0.17 ± 0.02	0.21 ± 0.02	0.42 ± 0.05
	ant-medium		0.14 ± 0.02	0.11 ± 0.01	0.17 ± 0.01	0.22 ± 0.05	0.09 ± 0.02	0.06 ± 0.01	0.54 ± 0.19
	ant-full-replay		0.18 ± 0.02	-0.07 ± 0.10	0.21 ± 0.02	0.20 ± 0.01	0.09 ± 0.01	0.11 ± 0.02	0.46 ± 0.10
	ant-expert		-0.16 ± 0.01	-0.23 ± 0.02	0.05 ± 0.04	0.16 ± 0.00	0.11 ± 0.03	0.10 ± 0.01	0.18 ± 0.00

The experimental procedure such as hyperparameters, training routine, and raw learning curves are detailed in Appendix B. All algorithms are trained to convergence and we used 10 seeds for the Iterated Coordination Game and 3 seeds for MAMuJoCo tasks. Tables are normalized and report the mean evaluation performance and the standard error of the mean across seeds. Evaluation is done for 100 episodes using the greedy policies (no sampling).

6.1 Strategy Agreement

Table 3 reports the results for the offline Iterated Coordination Game and validates most of our intuitions about strategy agreement: the centralized execution (IQL) and model-based (MOMA-PPO) approaches are able to coordinate agents regardless of the datasets. On the other hand, independent BC agents imitate the dataset behavior and therefore coordinate only if the dataset majorly demonstrates coordination. Surprisingly, the CTDE model-free approach MAIQL is able to break symmetry and coordinate agents in the neutral dataset. We hypothesize that small numerical errors in the centralized value approximation have the team favor one equivalent strategy over the other. Unfortunately, the conservatism of model-free methods forces agents to stay close to the demonstrated behaviors and prevails over this brittle symmetry-breaking mechanism in the unfavorable – i.e., uncoordinated – dataset.

Table 4 (a) confirms that these insights on strategy agreement hold in the more complex two-agent Reacher environment: model-free methods struggle with strategy agreement, especially under

partial observability¹. Again, CTDE methods – particularly MAIQL – tend to fare better than independent learners. Interestingly, IBC fares best among independent learners. Finally, our model-based CTDE approach *MOMA-PPO solves strategy agreement* and performs on par with centralized execution (IQL), a setting that sidesteps the strategy agreement issue altogether. MOMA-PPO matches or significantly outperforms all other baselines.

6.2 Strategy Fine-Tuning

From Table 4 (b) one can investigate how the different offline methods cope with strategy finetuning. First, IQL performs on par with the other model-free methods which suggests that centralized execution (single-agent) vs. decentralized execution (multi-agent) is less a consideration for strategy fine-tuning than it is for strategy agreement. Yet, this also highlights that *model-free methods (even when centralized) are unable to perform strategy fine-tuning*. Indeed, they are surpassed by our model-based method, MOMA-PPO. This latter generates additional synthetic experiences that *allow for strategy fine-tuning in addition to strategy agreement*. For model-free methods, independent learners tend to outperform CTDE ones (which echoes [43] observations). Additionally, comparing IQL with MAIQL performance does highlight that offline multi-agent coordination is more challenging in varied datasets (i.e. medium or full-replay that display both coordinated and uncoordinated behaviors) than in uniformly coordinated datasets (i.e. expert).

¹<https://sites.google.com/view/moma-ppo> shows independent learners converge to incompatible conventions.

In partially observable (PO) tasks, model-free methods are unable to adapt the behaviors demonstrated in the datasets and they result in teams that run in circles because the yellow agent (the only one to observe the torso’s headings and velocities) fails to correct the other limbs’ motions. Conversely, with MOMA-PPO, the yellow agent steers the ant toward the correct direction, and the teams reach very satisfactory performances provided that the datasets have enough coverage to learn a world model that can simulate diverse and robust behaviors (cf. the lower performance for the expert dataset).²

Finally, the poor performance of MAIQL-ft suggests that IQL’s finetuning abilities might not carry over to the multi-agent setting (even though we used the ground-truth simulator to generate the rollouts). While there might be multiple causes, we hypothesize that it is mainly due to MAIQL’s instability since it required intensive hyperparameters finetuning and filtering out collapsed runs for ant tasks. We believe that MAIQL’s instability is exacerbated by the induced non-stationarity of the training data when augmenting it with online interactions. Unlike online methods, offline algorithms are designed to learn on fixed datasets and are thus ill-equipped to deal with data continually collected by changing policies. Similarly, experiments that used MAIQL instead of MAPPO for MOMA (i.e. MOMA-IQL) quickly led to unstable learning and exploding losses.

In conclusion, our results show that current model-free offline MARL methods fail at offline coordination. Crucially, this deficiency remains even in very simple domains such as the Iterated Coordination Game. Also, it sometimes leads to underperforming basic imitation learning (i.e., IBC) or the datasets’ average score (see Table 2). Conversely, our model-based approach is able to coordinate agents even under severe partial observability and with learned world models by restoring inter-agent interactions throughout learning.

6.3 Ablations

We validate our design choice of using epistemic uncertainty over aleatoric. To do so we replace MOMA-PPO’s penalty with the aleatoric uncertainty penalty of MOPO [67]. For a fair comparison, we followed [67]’s parameter search procedure (i.e., coefficient values of 1 and 5). Additionally, we investigated aleatoric penalty both with and without the use of adaptive rollout length (which is based on epistemic uncertainty). The resulting learning curves on Ant-full-replay with full and partial observability are displayed in Figure 4. It appears that using epistemic uncertainty always significantly outperforms using aleatoric uncertainty except in the case of full observability and adaptive rollouts where the comparison is not statistically significant. This validates the soundness of our choice regarding the use of epistemic uncertainty over aleatoric uncertainty in deterministic environments.

Appendix B.2 reports additional ablations, which main insights we discuss here. First, using the ground-truth simulator yielded higher scores suggesting that the performance of MOMA-PPO can be further improved by learning a more accurate world model. Then, not clipping the generated next-states to the dataset’s bounding box prevented generating length 10 rollouts as the states’ magnitude

exploded after a few world-model steps. Also, the use of uncertainty penalty (λ_g) and adaptive rollouts’ length (l_ϵ) are necessary to achieve satisfactory results. Finally, varying the rollout’s maximum length from 5 to 50 did not significantly impact performance.

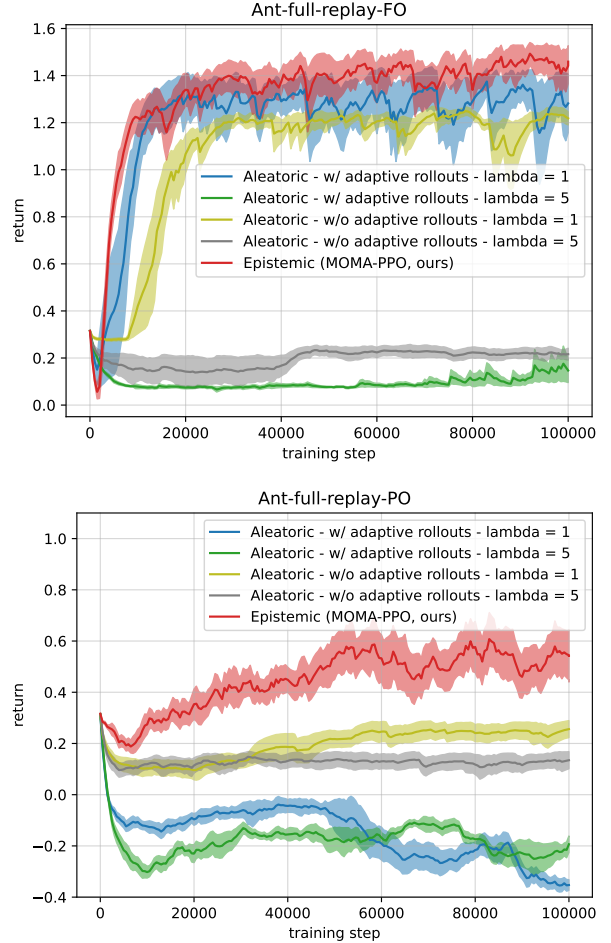


Figure 4: Comparison between using epistemic uncertainty reward penalty (MOMA-PPO) vs. aleatoric uncertainty reward penalty (MOPO-like). Mean and standard error of the mean on three seeds.

7 DISCUSSION AND CONCLUSION

This work explores coordination in offline MARL and highlights the failures of current model-free methods. For instance, they struggle in the presence of multiple equivalent but incompatible optimal team strategies (*strategy agreement*), or when partial observability requires the team to adapt the behaviors demonstrated in the dataset (*strategy fine-tuning*). In such scenarios, prevalent model-free methods might even fail to match the dataset’s performance and fall short compared to behavioral cloning. To address these problems, we propose MOMA-PPO which is, to our knowledge, the first model-based offline MARL approach. Our method is able to coordinate teams of offline learners and significantly outperforms model-free alternatives. Interestingly, it also outperforms the fully

²<https://sites.google.com/view/moma-ppo> shows rollouts with and without “steering” behavior.

centralized model-free method IQL [30] even though this latter completely bypasses the strategy agreement problem. This suggests that model-free methods, even when fully centralized, are unable to deal with strategy fine-tuning. We also observe in our experiments that the single-agent approach of fine-tuning offline methods (i.e. IQL) with online interactions [30] fails in the multi-agent setting (i.e. MAIQL). Our model-based approach succeeds at both multi-agent strategy agreement and strategy fine-tuning problems which goes to show that the benefits of offline model-based approaches over offline model-free ones [67] hold in the multi-agent setting. This might indicate that world models may generalize better than values networks in offline learning’s limited data regime. Finally, for tasks that require adapting the team’s behavior, dataset coverage might be more desirable than demonstrated performance. Indeed, methods fare better on random datasets than on expert ones.

MOMA-PPO’s successes put forward the benefits of model-based methods that leverage online policy optimization. Nevertheless, our approach of coordinating agents through a world model is general and some might be interested in extending it to more sample-efficient policy learning algorithms (i.e., MOMA-SAC). Also, the dataset bounding box clipping and adaptive rollouts developed for MOMA-PPO might benefit the single-agent setting.

Finally, this work aims to pinpoint an overlooked issue in the offline MARL community and proposes a new avenue of model-based solutions that shifts away from conventional model-free approaches. Therefore, we look forward to future work that will further analyze model-based offline MARL approaches and scale them to more domains and datasets.

ACKNOWLEDGMENTS

The authors would like to thank Luis Pineda, Hengyuan Hu, and Eugene Vinitsky for their help, insightful discussions, and advice. This research was enabled in part by support provided by Calcul Québec and the Digital Research Alliance of Canada. This work was partially conducted while Paul was interning at Meta AI - FAIR.

REFERENCES

- [1] Dzmitry Bahdanau, Kyung Hyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. In *3rd International Conference on Learning Representations, ICLR 2015*.
- [2] Paul Barde, Tristan Karch, Derek Nowrouzezahrai, Clément Moulin-Frier, Christopher Pal, and Pierre-Yves Oudeyer. 2022. Learning to Guide and to be Guided in the Architect-Builder Problem. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=swiyAeGzFhQ>
- [3] Paul Barde, Julien Roy, Wonseok Jeon, Joelle Pineau, Chris Pal, and Derek Nowrouzezahrai. 2020. Adversarial soft advantage fitting: Imitation learning without policy optimization. *Advances in Neural Information Processing Systems* 33 (2020), 12334–12344.
- [4] Samuel Barrett, Peter Stone, and Sarit Kraus. 2011. Empirical evaluation of ad hoc teamwork in the pursuit domain. In *The 10th International Conference on Autonomous Agents and Multiagent Systems-Volume 2*, 567–574.
- [5] Craig Boutilier. 1996. Planning, learning and coordination in multiagent decision processes. In *TARK*, Vol. 96. Citeseer, 195–210.
- [6] Justin Boyan and Michael Littman. 1993. Packet routing in dynamically changing networks: A reinforcement learning approach. *Advances in neural information processing systems* 6 (1993).
- [7] Ronen Brafman and Moshe Tennenholtz. 2002. Efficient learning equilibrium. *Advances in Neural Information Processing Systems* 15 (2002).
- [8] Wilfried Brauer and Gerhard Weiß. 1998. Multi-machine scheduling—a multi-agent learning approach. In *Proceedings International Conference on Multi Agent Systems (Cat. No. 98EX160)*. IEEE, 42–48.
- [9] Micah Carroll, Rohin Shah, Mark K Ho, Tom Griffiths, Sanjit Seshia, Pieter Abbeel, and Anca Dragan. 2019. On the utility of learning about humans for human-ai coordination. *Advances in neural information processing systems* 32 (2019).
- [10] Lars-Erik Cederman. 1997. *Emergent actors in world politics: how states and nations develop and dissolve*. Vol. 2. Princeton University Press.
- [11] Georgios Chalkiadakis and Craig Boutilier. 2003. Coordination in multiagent reinforcement learning: A Bayesian approach. In *Proceedings of the second international joint conference on Autonomous agents and multiagent systems*, 709–716.
- [12] Caroline Claus and Craig Boutilier. 1998. The dynamics of reinforcement learning in cooperative multiagent systems. *AAAI/IAAI 1998*, 746–752 (1998), 2.
- [13] Vincent P Crawford and Joel Sobel. 1982. Strategic information transmission. *Econometrica: Journal of the Econometric Society* (1982), 1431–1451.
- [14] Christian Schroeder de Witt, Tarun Gupta, Denys Makoviichuk, Viktor Makoviychuk, Philip HS Torr, Mingfei Sun, and Shimon Whiteson. 2020. Is independent learning all you need in the starcraft multi-agent challenge? *arXiv preprint arXiv:2011.09533* (2020).
- [15] Kurt Dresner and Peter Stone. 2004. Multiagent traffic management: A reservation-based intersection control mechanism. In *Autonomous Agents and Multiagent Systems, International Joint Conference on*, Vol. 3. Citeseer, 530–537.
- [16] K Fischer, N Kuhn, HJ Muller, JP Muller, and M Pischel. 1993. Sophisticated and distributed: The transportation domain. In *Proceedings of 9th IEEE Conference on Artificial Intelligence for Applications*. IEEE, 454.
- [17] Jakob Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson. 2018. Counterfactual multi-agent policy gradients. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 32.
- [18] Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. 2020. D4rl: Datasets for deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219* (2020).
- [19] Scott Fujimoto and Shixiang Shane Gu. 2021. A minimalist approach to offline reinforcement learning. *Advances in neural information processing systems* 34 (2021), 20132–20145.
- [20] Scott Fujimoto, David Meger, and Doina Precup. 2019. Off-policy deep reinforcement learning without exploration. In *International conference on machine learning*. PMLR, 2052–2062.
- [21] Stephen Grand, Dave Cliff, and Anil Malhotra. 1997. Creatures: Artificial life autonomous software agents for home entertainment. In *Proceedings of the first international conference on Autonomous agents*, 22–29.
- [22] David Ha and Jürgen Schmidhuber. 2018. World models. *arXiv preprint arXiv:1803.10122* (2018).
- [23] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. 2018. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*. PMLR, 1861–1870.
- [24] Hengyuan Hu, Adam Lerer, Alex Peysakhovich, and Jakob Foerster. 2020. “Other-Play” for Zero-Shot Coordination. In *International Conference on Machine Learning*. PMLR, 4399–4410.
- [25] Jun Huang, Nicholas R Jennings, and John Fox. 1995. An agent architecture for distributed medical care. In *Intelligent Agents: ECAI-94 Workshop on Agent Theories, Architectures, and Languages Amsterdam, The Netherlands August 8–9, 1994 Proceedings 1*. Springer, 219–232.
- [26] Natasha Jaques, Angeliki Lazaridou, Edward Hughes, Caglar Gulcehre, Pedro Ortega, Dj Strouse, Joel Z Leibo, and Nando De Freitas. 2019. Social Influence as Intrinsic Motivation for Multi-Agent Deep Reinforcement Learning. In *International Conference on Machine Learning*, 3040–3049.
- [27] Jiechuan Jiang and Zongqing Lu. 2021. Offline decentralized multi-agent reinforcement learning. *arXiv preprint arXiv:2108.01832* (2021).
- [28] Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims. 2020. Morel: Model-based offline reinforcement learning. *Advances in neural information processing systems* 33 (2020), 21810–21823.
- [29] Diederik P Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. *International Conference on Learning Representations* (2015).
- [30] Ilya Kostrikov, Ashvin Nair, and Sergey Levine. 2021. Offline Reinforcement Learning with Implicit Q-Learning. *arxiv* (2021).
- [31] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. 2020. Conservative q-learning for offline reinforcement learning. *Advances in Neural Information Processing Systems* 33 (2020), 1179–1191.
- [32] Angeliki Lazaridou, Alexander Peysakhovich, and Marco Baroni. 2017. Multi-Agent Cooperation and the Emergence of (Natural) Language. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=Hk8N3ScIq>
- [33] Adam Lerer and Alexander Peysakhovich. 2019. Learning existing social conventions via observationally augmented self-play. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, 107–114.
- [34] Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. 2020. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643* (2020).
- [35] Jiaoyang Li, Zhe Chen, Yi Zheng, Shao-Hung Chan, Daniel Harabor, Peter J Stuckey, Hang Ma, and Sven Koenig. 2021. Scalable rail planning and replanning: Winning the 2020 flatland challenge. In *Proceedings of the International Conference on Automated Planning and Scheduling*, Vol. 31, 477–485.

- [36] Michael L Littman et al. 2001. Friend-or-foe Q-learning in general-sum games. In *ICML*, Vol. 1. 322–328.
- [37] Haotian Liu and Wenchuan Wu. 2021. Online multi-agent reinforcement learning for decentralized inverter-based volt-var control. *IEEE Transactions on Smart Grid* 12, 4 (2021), 2980–2990.
- [38] Ryan Lowe, Yi Wu, Aviv Tamar, Jean Harb, Pieter Abbeel, and Igor Mordatch. 2017. Multi-agent actor-critic for mixed cooperative-competitive environments. In *Advances in Neural Information Processing Systems*. 6379–6390.
- [39] Xueguang Lyu, Yuchen Xiao, Brett Daley, and Christopher Amato. 2021. Contrasting centralized and decentralized critics in multi-agent reinforcement learning. *arXiv preprint arXiv:2102.04402* (2021).
- [40] Anuj Mahajan, Tabish Rashid, Mikayel Samvelyan, and Shimon Whiteson. 2019. Maven: Multi-agent variational exploration. *Advances in Neural Information Processing Systems* 32 (2019).
- [41] Ranjit Nair, Milind Tambe, Makoto Yokoo, David Pynadath, and Stacy Marsella. 2003. Taming decentralized POMDPs: Towards efficient policy computation for multiagent settings. In *IJCAI*, Vol. 3. Citeseer, 705–711.
- [42] Praveen Palanisamy. 2020. Multi-agent connected autonomous driving using deep reinforcement learning. In *2020 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 1–7.
- [43] Ling Pan, Longbo Huang, Tengyu Ma, and Huazhe Xu. 2022. Plan better amid conservatism: Offline multi-agent reinforcement learning with actor rectification. In *International Conference on Machine Learning*. PMLR, 17221–17237.
- [44] Liviu Panait and Sean Luke. 2005. Cooperative multi-agent learning: The state of the art. *Autonomous agents and multi-agent systems* 11, 3 (2005), 387–434.
- [45] Bei Peng, Tabish Rashid, Christian Schroeder de Witt, Pierre-Alexandre Kamieny, Philip Torr, Wendelin Böhrer, and Shimon Whiteson. 2021. Facmac: Factored multi-agent centralized policy gradients. *Advances in Neural Information Processing Systems* 34 (2021), 12208–12221.
- [46] Dean A Pomerleau. 1991. Efficient training of artificial neural networks for autonomous navigation. *Neural computation* 3, 1 (1991), 88–97.
- [47] Tabish Rashid, Mikayel Samvelyan, Christian Schroeder, Gregory Farquhar, Jakob Foerster, and Shimon Whiteson. 2018. Qmix: Monotonic value function factorisation for deep multi-agent reinforcement learning. In *International Conference on Machine Learning*. PMLR, 4295–4304.
- [48] Julien Roy, Paul Barde, Félix Harvey, Derek Nowrouzezahrai, and Chris Pal. 2020. Promoting coordination through policy regularization in multi-agent deep reinforcement learning. *Advances in Neural Information Processing Systems* 33 (2020), 15774–15785.
- [49] John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. 2015. High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438* (2015).
- [50] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347* (2017).
- [51] Tianyu Shi, Dong Chen, Kaian Chen, and Zhaojian Li. 2021. Offline Reinforcement Learning for Autonomous Driving with Safety and Exploration Enhancement. *arXiv preprint arXiv:2110.07067* (2021).
- [52] Jonathan Spencer, Sanjiban Choudhury, Arun Venkatraman, Brian Ziebart, and J Andrew Bagnell. 2021. Feedback in imitation learning: The three regimes of covariate shift. *arXiv preprint arXiv:2102.02872* (2021).
- [53] Randall Steeb, Stephanie Cammarata, Frederick A Hayes-Roth, Perry W Thorndyke, and Robert E Wesson. 1981. *Distributed intelligence for air fleet control*. Technical Report. RAND CORP SANTA MONICA CA.
- [54] Peter Stone, Gal Kaminka, Sarit Kraus, and Jeffrey Rosenschein. 2010. Ad hoc autonomous agent teams: Collaboration without pre-coordination. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 24. 1504–1509.
- [55] Richard S Sutton. 1990. Integrated architectures for learning, planning, and reacting based on approximating dynamic programming. In *Machine learning proceedings 1990*. Elsevier, 216–224.
- [56] Wei-Cheng Tseng, Tsun-Hsuan Johnson Wang, Yen-Chen Lin, and Phillip Isola. 2022. Offline Multi-Agent Reinforcement Learning with Knowledge Distillation. *Advances in Neural Information Processing Systems* 35 (2022), 226–237.
- [57] Mycal Tucker, Yilun Zhou, and Julie Shah. 2020. Adversarially guided self-play for adopting social conventions. *arXiv preprint arXiv:2001.05994* (2020).
- [58] László Varga, Nick R Jennings, and David Cockburn. 1994. Integrating intelligent systems into a cooperating community for electricity distribution management. *Expert Systems with Applications* 7, 4 (1994), 563–579.
- [59] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).
- [60] Eugene Vinitzky, Nathan Lichtlé, Xiaomeng Yang, Brandon Amos, and Jakob Foerster. 2022. Nocturne: a scalable driving benchmark for bringing multi-agent learning one step closer to the real world. *arXiv preprint arXiv:2206.09889* (2022).
- [61] Jianhao Wang, Wenzhe Li, Haozhe Jiang, Guangxiang Zhu, Siyuan Li, and Chongjie Zhang. 2021. Offline reinforcement learning with reverse model-based imagination. *Advances in Neural Information Processing Systems* 34 (2021), 29420–29432.
- [62] Daniël Willemsen, Mario Coppola, and Guido CHE de Croon. 2021. MAMBPO: Sample-efficient multi-robot reinforcement learning using learned world models. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 5635–5640.
- [63] Mark Woodward, Chelsea Finn, and Karol Hausman. 2020. Learning to interactively learn and assist. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 34. 2535–2543.
- [64] Yiqin Yang, Xiaoteng Ma, Chenghao Li, Zewu Zheng, Qiyuan Zhang, Gao Huang, Jun Yang, and Qianchuan Zhao. 2021. Believe What You See: Implicit Constraint Approach for Offline Multi-Agent Reinforcement Learning. *NeurIPS* (2021).
- [65] Chao Yu, Akash Velu, Eugene Vinitzky, Jiaxuan Gao, Yu Wang, Alexandre Bayen, and Yi Wu. 2022. The surprising effectiveness of ppo in cooperative multi-agent games. *Advances in Neural Information Processing Systems* 35 (2022), 24611–24624.
- [66] Tianhe Yu, Aviral Kumar, Rafael Rafailov, Aravind Rajeswaran, Sergey Levine, and Chelsea Finn. 2021. Combo: Conservative offline model-based policy optimization. *Advances in neural information processing systems* 34 (2021), 28954–28967.
- [67] Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Y Zou, Sergey Levine, Chelsea Finn, and Tengyu Ma. 2020. Mopo: Model-based offline policy optimization. *Advances in Neural Information Processing Systems* 33 (2020), 14129–14142.
- [68] Chi Zhang, Sanmukh R Kuppannagari, Chuanxiu Xiong, Rajgopal Kannan, and Viktor K Prasanna. 2019. A cooperative multi-agent deep reinforcement learning framework for real-time residential load scheduling. In *Proceedings of the International Conference on Internet of Things Design and Implementation*. 59–69.
- [69] Chongjie Zhang and Victor Lesser. 2013. Coordinating multi-agent reinforcement learning with limited communication. In *Proceedings of the 2013 international conference on Autonomous agents and multi-agent systems*. 1101–1108.
- [70] Qizhen Zhang, Chris Lu, Animesh Garg, and Jakob Foerster. 2022. Centralized Model and Exploration Policy for Multi-Agent RL. In *Proceedings of the 21st International Conference on Autonomous Agents and Multiagent Systems (Virtual Event, New Zealand) (AAMAS '22)*. International Foundation for Autonomous Agents and Multiagent Systems, Richland, SC, 1500–1508.
- [71] Weinan Zhang, Xihuai Wang, Jian Shen, and Ming Zhou. 2021. Model-based Multi-agent Policy Optimization with Adaptive Opponent-wise Rollouts. *International Joint Conference on Artificial Intelligence* (2021).

A REPRODUCIBILITY DETAILS

The following section focuses on reproducibility and goes into detail about the implementations and experimental procedures.

A.1 Methods implementation

MOMA-PPO entropy bonus and action penalty. For offline methods based on online RL algorithm, exploration is an important component (cf. TD3’s exploration strategy) so we used an entropy bonus for PPO defined as:

```
# entropy bonus
# dimensions are batch, act_dim, n_agents.
# Instead of using closed form entropy,
# we estimate it with E(-pi log pi) were
# the expectation is sampled over pi_old
# so we have to correct with pi_new/pi_old (which is ratio!)

# we do two losses
# (clipped / not clipped like with the actor loss)
surrogate_entropy = - (ratio * new_policy).mean(0)
clipped_entropy = - (clipped_ratio * new_policy).mean(0)

entropy = torch.min(surrogate_entropy, clipped_entropy)

self.entropy_alpha += self.ppo_entropy_bonus_coeff*(self.ppo_entropy_target - entropy).detach()
self.entropy_alpha.data.clamp_min_(0.)

# minus sign because we minimize the expressions
# i.e. max(ent) = min(-ent)
entropy_bonus = - entropy * self.entropy_alpha
```

with an entropy bonus coefficient of 0.001 and an entropy target of -6 and -4 for respectively Ant and Reacher tasks. Since entropy computation can become numerically unstable for squashed actions close to the Tanh bounds, PPO uses action penalty instead of Tanh squashing to keep actions close to the -1, 1 range:

```
delta = (1. - actions.abs())
action_bound_error = ((delta < 0.).to(torch.float32) * delta**2).sum(1, keepdim=True)

surrogate_action_bound_error = (ratio * action_bound_error).mean(0)
clipped_action_bound_error = (clipped_ratio * action_bound_error).mean(0)

action_bound_error = torch.max(surrogate_action_bound_error, clipped_action_bound_error)
action_penalty = self.ppo_action_penalty_coeff * action_bound_error
```

with an action penalty coefficient of 1.

General Advantage Estimation. We show below how we modified the general advantage estimation to account for rollout termination (indicated by `time_out_masks`).

```
# initial (end) running returns is the next state value
running_returns = next_values[-1] * masks[-1]

# initial (end) advantage is 0 because no difference in value and return
running_advants = 0

for t in reversed(range(0, len(rewards))):

    # We are going reverse so if done, only reward because end of
    # episode if timeout, we stop accumulation and use value as
    # bootstrap like in initialization
    running_returns = rewards[t] + self.discount * masks[t] * (running_returns * time_out_masks[t] + (1. -
        time_out_masks[t]) * next_values[t] * masks[t])

    returns[t] = running_returns

    ## No accumulation here and timeout doesn't influence next_state value
    running_delta = rewards[t] + (self.discount * next_values[t] * masks[t]) - values[t]

    ## if timeout running_advants goes back to running_delta because
    # we do not have extra rewards to estimate it
    # (cf initialization above)
    running_advants = running_delta + (self.discount * self.lamda * running_advants * masks[t]) * time_out_masks[t]

    advants[t] = running_advants
```

Memory module for partial observability. To handle partial observability we use observation-action histories h_t^i of size ten (i.e. the ten past observation-action pairs). These histories are processed with self-attention followed by soft-attention to yield embeddings of size $e_h = 128$ that are concatenated to the current state before being fed to the policy and value networks. We use a first linear layer to encode h_t^i to \mathbb{R}^{e_h} and add positional encodings [59]. Query, Key, and Value networks are linear layers and we follow [59]’s Scaled Dot-Product Attention with skip connection and layer-norm. Finally, the resulting self-attentions are aggregated using soft-attention with the soft-key network being a bias-less linear layer and the soft-queries are e_h normally initialized trainable parameters.

Note that policy and value networks use (and backprop through) the same memory module but target networks have their own target memory module (that tracks the memory module with Polyak updates just like regular target networks). The memory learning rate is $1e^{-4}$ for all the algorithms.

MAIQL. We provide IQL [30] learning rules here for completeness. Value learning is done with SARSA Bellman on e expectile of Q instead of mean Q (this latter is a special case where $e = 0.5$). To ensure that the expectile is computed only from the action selection distribution and is not influenced by the randomness of the environment’s transitions, IQL uses a state-only value function V that marginalizes over future transitions:

$$\begin{aligned} L_V(\psi) &= \mathbb{E}_{s,a \sim \mathcal{D}} \left[\mathcal{L}_2^e(Q_{\hat{\theta}}(s, a) - V_{\psi}(s)) \right] \\ &= \mathbb{E}_{s,a \sim \mathcal{D}} \left[|e - \mathbb{I}(Q_{\hat{\theta}}(s, a) - V_{\psi}(s) < 0)| (Q_{\hat{\theta}}(s, a) - V_{\psi}(s))^2 \right], \end{aligned} \quad (1)$$

$$L_Q(\theta) = \mathbb{E}_{s,a,r,s' \sim \mathcal{D}} \left[\left(r(s, a) + \gamma V_{\psi}(s') - Q_{\theta}(s, a) \right)^2 \right]. \quad (2)$$

Policy extraction is done with Advantage Weighted Regression (AWR):

$$L_{\pi}(\phi) = \mathbb{E}_{(s,a) \sim \mathcal{D}} \left[-\exp\left(\beta(Q_{\hat{\theta}}(s, a) - V_{\psi}(s))\right) \log \pi_{\phi}(a|s) \right]. \quad (3)$$

We make IQL multi-agent (i.e. MAIQL) by leveraging the CTDE formalism and using QMIX value-decomposition [47] for both Q and V :

$$\begin{aligned} V_{\psi}(s) &= \sum_i w_V^i(s) V_{\psi^i}(h^i) + b_V(s), \\ Q_{\theta}(s) &= \sum_i w_Q^i(s) Q_{\theta^i}(h^i) + b_Q(s). \end{aligned} \quad (4)$$

And the target network:

$$Q_{\hat{\theta}}(s) = \sum_i \hat{w}_Q^i(s) Q_{\hat{\theta}^i}(s^i) + \hat{b}_Q(s). \quad (5)$$

Similarly, the joint policy is assumed to decompose as:

$$\pi(a|s) \triangleq \prod_i \pi_i(a^i|h^i). \quad (6)$$

Injecting this into Eq. 3 one gets:

$$\begin{aligned} L_{\pi}(\phi) &= \mathbb{E}_{(s,a) \sim \mathcal{D}} \left[-B(s) \prod_i \exp\left(\beta(\hat{w}_Q^i(s) Q_{\hat{\theta}^i}(h^i, a^i) - w_V^i(s) V_{\psi^i}(h^i))\right) \sum_j \log \pi_{\phi^j}(a^j|h^j) \right], \\ B(s) &= \exp(\hat{b}_Q(s) - b_V(s)). \end{aligned} \quad (7)$$

Finetuning MAIQL: MAIQL-ft. We follow [30]’s finetuning procedure and use the ground truth simulator to generate the rollout. However, the rollouts are still generated using MOMA’s Dyna-like approach described in Section 3 rather than generating length 1000 rollouts from the initial state distribution. Indeed, we consider the model-based offline setting and not the offline-to-online finetuning setting. Therefore it is unfeasible to assume that the task’s ground-truth initial state distribution is known or that it is possible to learn a world model that remains accurate over 1000 simulation steps.

Opensource baselines. For the baseline implementations, we followed the official repositories at:

- https://github.com/sfujim/TD3_BC,
- <https://github.com/ling-pan/OMAR/>,
- and https://github.com/ikostrikov/implicit_q_learning.

A.2 Hyperparameters, Tuning, and Training

Unless specified otherwise, networks are two-layers 256 units ReLu MLPs. We use Adam optimizer [29] with default hyperparameters (except learning rates) and a batch size of 256. We clip gradient norms to 1.

World model learning. World models use four layers and 1024 units. World models use a learning rate of $3e^{-5}$ and are trained for $3e^6$ steps.

Policy learning. MAPPO learning rate was finetuned on MAMUJOCO online task halfcheetah-v2_2x3_full with a grid search across $[1e^{-6}, 5e^{-6}, 1e^{-5}, 5e^{-5}, 1e^{-4}, 5e^{-4}, 1e^{-3}]$. We kept the value of $5e^{-5}$ for MOMA-PPO for all experiments. PPO uses rollouts length of 1000, 5 epochs per update, 2000 transitions between updates, a clip value of 0.2, a λ value of 0.98, and a critic loss coefficient of 0.5.

MAIQL showed quite unstable and we had to finetune its learning rate extensively depending on the tasks. We tried the range $[1e^{-4}, 3e^{-4}, 1e^{-5}, 5e^{-5}, 1e^{-6}]$ on ant-expert for MAIQL and MAIQL-ft and kept $3e^{-4}$. For the results on ant-expert with full observability, we had to discard one collapsed unstable seed for MAIQL and retrain another seed. For reacher tasks we tried $[5e^{-5}, 1e^{-4}, 3e^{-4}, 5e^{-4}, 1e^{-3}]$ and kept $3e^{-4}$. IQL and MAIQL use an expectile value of 0.7 and an AWR temperature of 3.

All other methods use their default learning rate of $3e^{-4}$. ITD3+BC uses a BC regularization parameter of 2.5, a policy frequency update of 2, a policy noise of 0.2, and a noise clip value of 0.5. ICQL uses a coefficient of 1 (and so does IOMAR’s CQL component) for all tasks except for ant-expert tasks where we had to tune it between $[0.1, 0.5, 1, 5]$ and kept a value of 5 (same for IOMAR). Additionally, CQL uses an LSE temperature of 1, 10 sampled actions, and a sample noise level of 0.2. Finally, IOMAR uses a coefficient of 1, 2 iterations, a μ value of 0, a σ value of 2, 20 samples, and 5 elites.

We use default values for other hyperparameters and we report them here for consistency. We use a discount factor of 0.99 and a target update coefficient of 0.005 for the Polyak averaging. TD3 (and thus CQL and OMAR) policies are deterministic with Tanh squashing while IQL and PPO use Gaussians with state-dependent variance. IQL uses Tanh squashing while PPO does not (see below).

Model-free methods are trained for $1e^6$ learning steps while MOMA-PPO is trained for $1e^5$ learning steps which correspond to roughly $2e^8$ collected interactions from the world model.

A.3 Compute

Our longest MOMA-PPO training took 6 days on a Tesla P100-PCIE-12GB GPU. For comparison our longest IOMAR run took 38 hours on a Tesla V100-SXM2-32GB GPU and our longest MAIQL run took 40 hours on a Tesla V100-SXM2-32GB GPU and five days on a Tesla P100-PCIE-12GB GPU for MAIQL-ft. Training a world model took at most three days on a Tesla P100-PCIE-12GB GPU. Note that training times are also impacted by how often we estimate performance for the learning curves and we do it much more often for MOMA-PPO (5 episodes every 50 learning steps) than for other methods (10 episodes every 5000 learning steps).

B RAW RESULTS

This section focuses on transparency and provides the raw results of the experiments.

B.1 Learning Curves

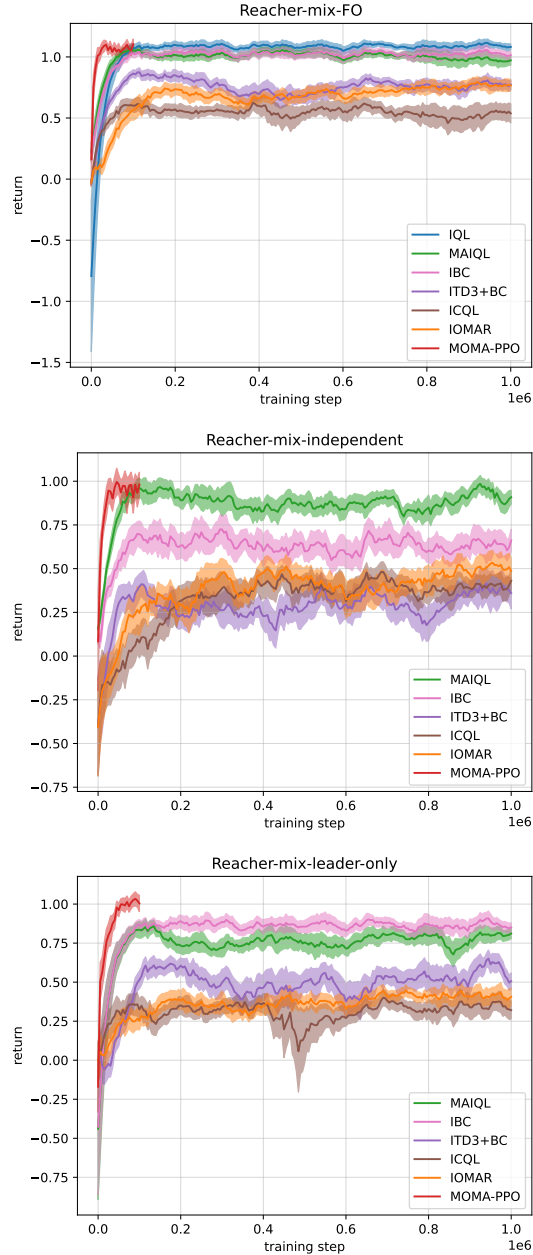


Figure 5: Learning Curves for two-agent Reacher. Mean and standard error of the mean on three seeds.

Figures 5 and 6 show the learning curves for the Reacher and Ant environments respectively. We use 5 episodes to evaluate MOMA-PPO every 50 learning steps and 10 episodes every 5000 learning steps to evaluate the other methods (this is why MOMA-PPO curves might look noisier). We used a smoothing factor of .8 for all curves. We report the mean over three seeds and the shaded area represents \pm the standard error of the mean. We train MOMA-PPO for $1e^5$ training steps because it is on-policy while the off-policy methods are trained for $1e^6$ training steps. MAIQL-ft training is $2e^6$ steps (half offline and half online fine-tuning).

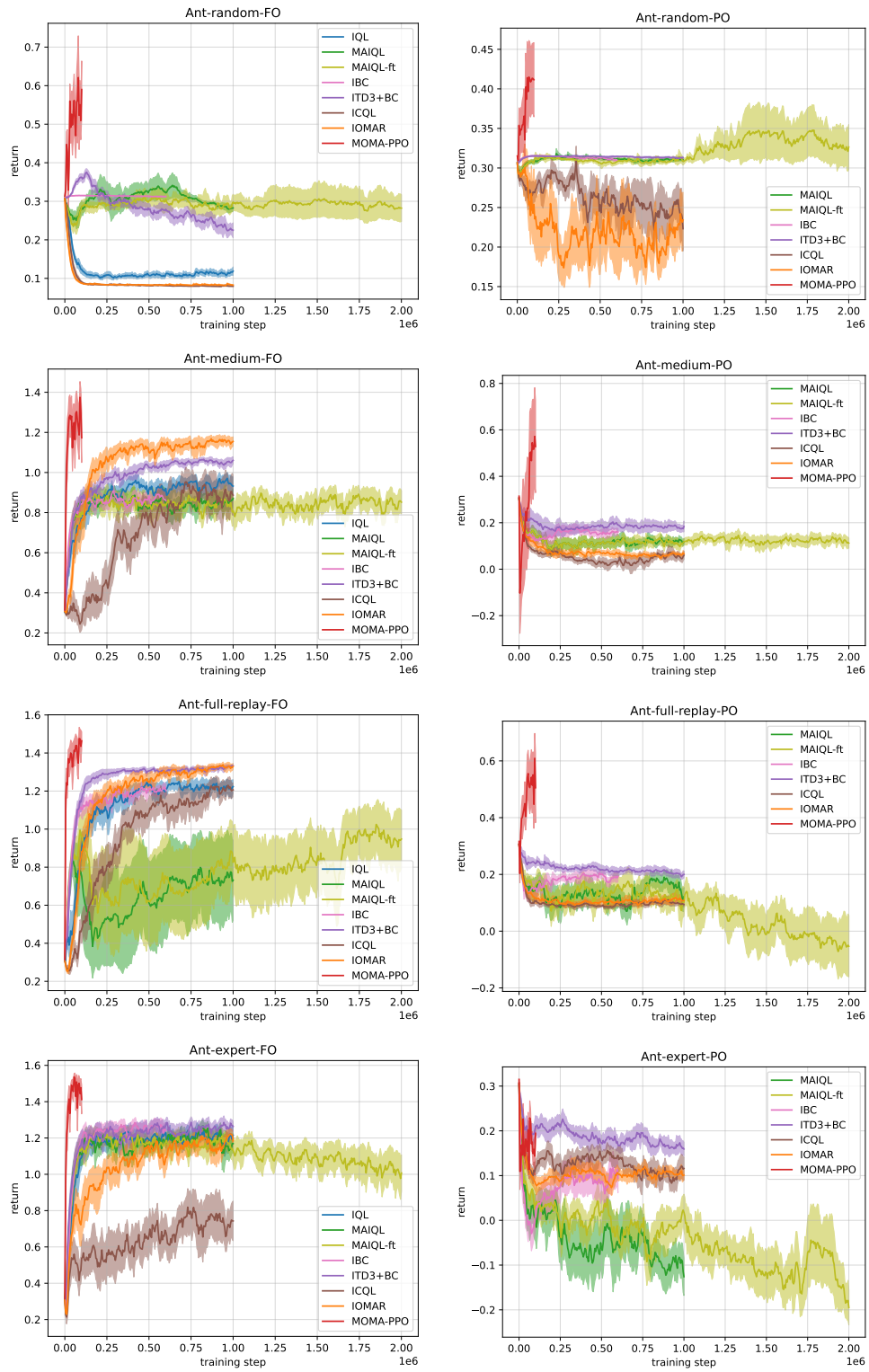


Figure 6: Learning curves for four-agent Ant. Mean and standard error of the mean on three seeds.

B.2 Ablations

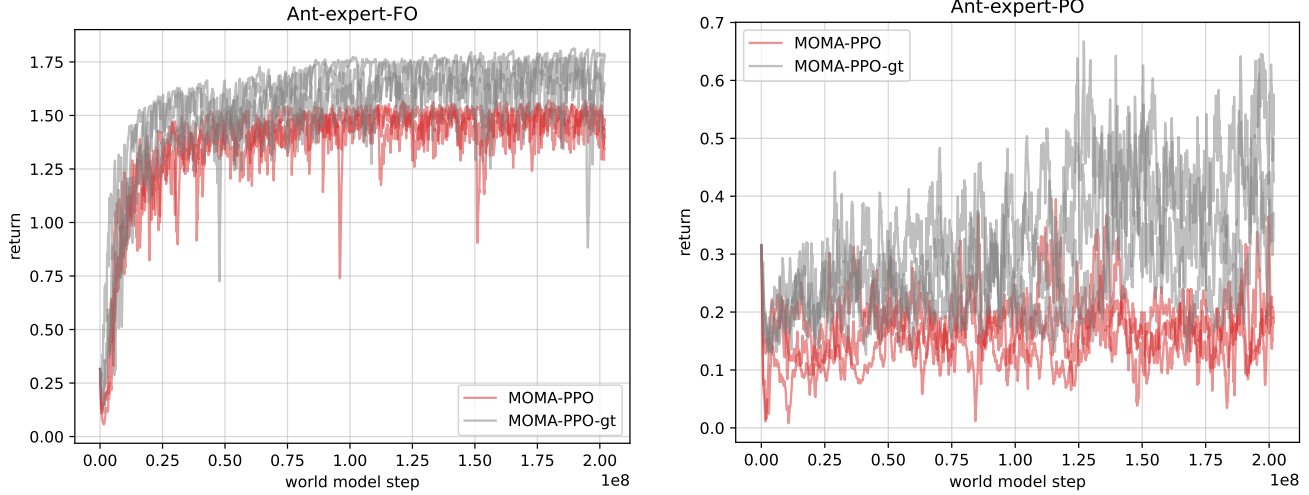


Figure 7: Impact of using the ground truth world model vs. the learned world model. “gt” stands for “ground truth” and means that the corresponding runs use the ground truth world model. The learning curves are with respect to generated transitions (either from the ground truth or the learned world model). Here we show each individual run instead of the usual mean and standard error of the mean.

Table 5: Mean scores and standard error of the mean at the end of training with and without the use of the ground truth simulator. Evaluations are over 100 episodes.

		MOMA-PPO	MOMA-PPO-gt
(FO)	ant-random	0.52 ± 0.07	1.17 ± 0.03
	ant-medium	1.29 ± 0.06	1.68 ± 0.03
	ant-full-replay	1.42 ± 0.07	1.66 ± 0.03
	ant-expert	1.49 ± 0.01	1.71 ± 0.04
(PO)	ant-random	0.42 ± 0.05	0.47 ± 0.01
	ant-medium	0.54 ± 0.19	0.81 ± 0.20
	ant-full-replay	0.46 ± 0.10	0.84 ± 0.07
	ant-expert	0.18 ± 0.00	0.43 ± 0.06

Figure 7 and Table 5 compare using a learned world model with having access to the ground truth simulator to generate the rollouts. It appears that MOMA-PPO performance can be further improved provided that we learn better world models.

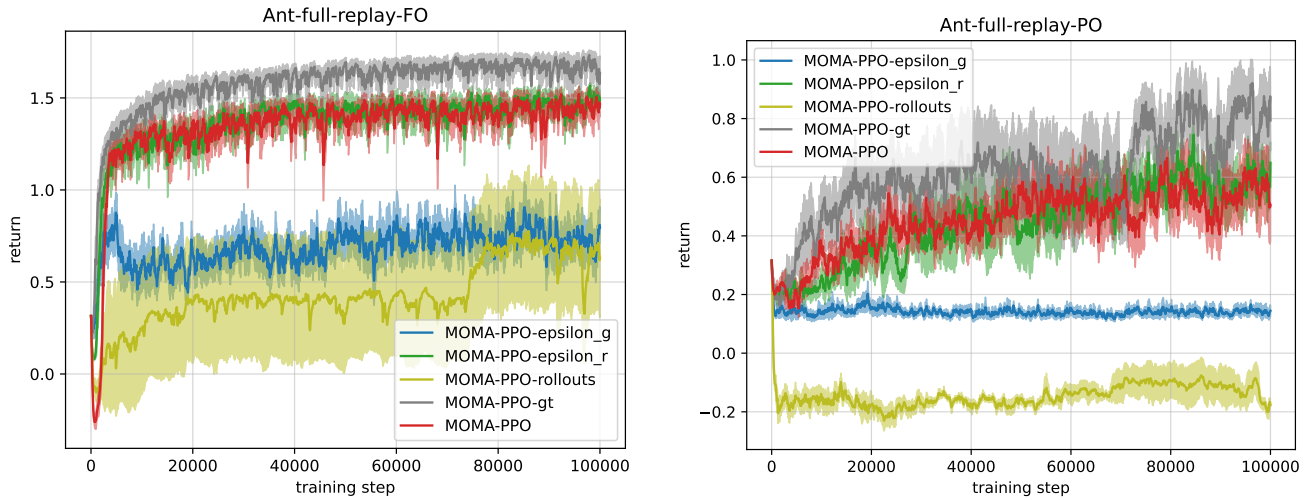


Figure 8: Impact of the different uncertainty-based techniques to prevent model exploitation by the learning algorithm. “gt” stands for “ground truth” and means that the corresponding runs use the ground truth world model. “epsilon_g” and “epsilon_r” respectively means that we set $\lambda_g = 0$ and $\lambda_r = 0$. “rollouts” means that we removed the adaptive rollouts component (i.e. $l_e = \infty$). Mean and standard error of the mean on three seeds.

Figure 8 shows ablations on MOMA-PPO. It appears that λ_r penalty can be removed without hurting performance, indeed it is already encapsulated in λ_g . The other components of MOMA-PPO such as λ_g uncertainty penalty on the reward, or the adaptive rollouts that terminate if the uncertainty crosses a threshold, are mandatory to reach satisfactory performance. Note that running experiments without clipping the world model predictions to the datasets’ bounding box was impossible as the magnitude of the predicted state exploded after a few rollout steps. This could be mitigated by reducing the exploration of the PPO policies so that the agents stay closer to the dataset distribution where the world model does not predict such extreme states.

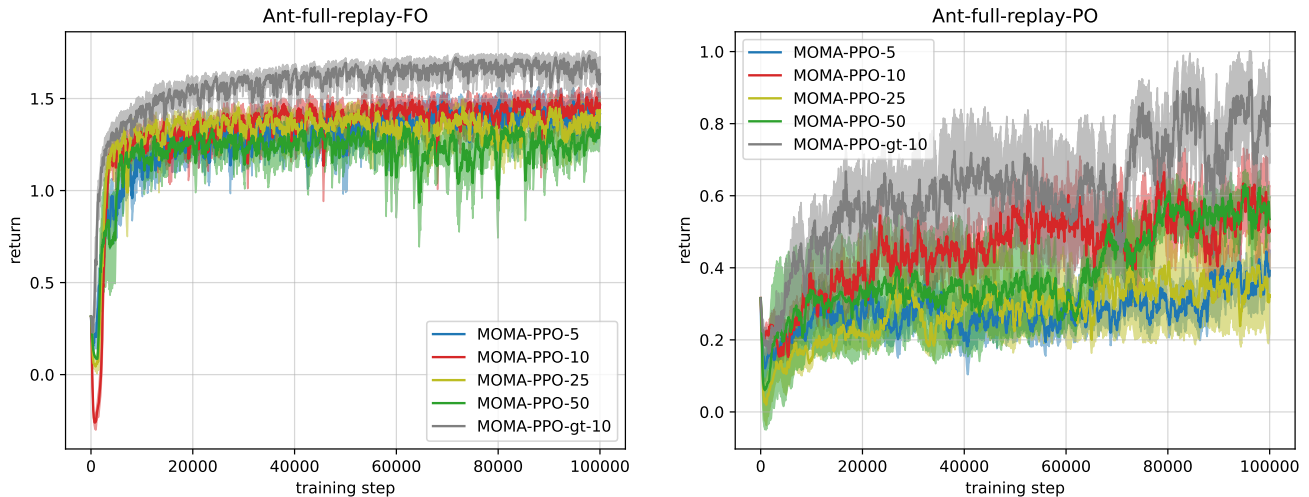


Figure 9: Impact of different maximum rollout lengths when generating the synthetic interactions. MOMA-PPO-x indicates a maximum rollout length of x. “gt” stands for “ground truth” and means that the corresponding runs use the ground truth world model. Mean and standard error of the mean on three seeds.

Figure 9 shows MOMA-PPO trainings for different maximum rollout lengths. It appears that varying the maximum length from 5 to 50 does not significantly impact MOMA-PPO performance on the tasks we investigated.