DECENTRALIZED TRANSFORMERS WITH CENTRALIZED AGGREGATION ARE SAMPLE-EFFICIENT MULTI-AGENT WORLD MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Learning a world model for model-free Reinforcement Learning (RL) agents can significantly improve the sample efficiency by learning policies in imagination. However, building a world model for Multi-Agent RL (MARL) can be particularly challenging due to the scalability issue in a centralized architecture arising from a large number of agents, and also the non-stationarity issue in a decentralized architecture stemming from the inter-dependency among agents. To address both challenges, we propose a novel world model for MARL that learns decentralized local dynamics for scalability, combined with a centralized representation aggregation from all agents. We cast the dynamics learning as an auto-regressive sequence modeling problem over discrete tokens by leveraging the expressive Transformer architecture, in order to model complex local dynamics across different agents and provide accurate and consistent long-term imaginations. As the first pioneering Transformer-based world model for multi-agent systems, we introduce a Perceiver Transformer as an effective solution to enable centralized representation aggregation within this context. Main results on Starcraft Multi-Agent Challenge (SMAC) and additional results on MAMujoco show that it outperforms strong model-free approaches and existing model-based methods in both sample efficiency and overall performance.

032

006

008 009 010

011

013

014

015

016

017

018

019

021

025

026

027

1 INTRODUCTION

Multi-Agent Reinforcement Learning (MARL) has made remarkable progress, which was driven largely by model-free algorithms (Nguyen et al., 2020). However, due to the complexity of multi-agent systems arising from large state-action space and partial observability, such algorithms usually demand extensive interactions to learn coordinative behaviors (Hernandez-Leal et al., 2020). A promising solution is building a world model that approximates the environment, which has exhibited its superior sample efficiency compared to model-free approaches in single-agent RL (Hafner et al., 2020; Łukasz Kaiser et al., 2020; Hafner et al., 2021; 2023; Hansen et al., 2022; 2024). However, extending the design of world model in single-agent domain to the multi-agent context encounters significant challenges due to the unique biases and characteristics inherent to multi-agent environments.

The challenges primarily stem from two different means for multi-agent dynamics learning: central-042 *ized* and *decentralized*. Learning a world model to approximate the *centralized* dynamics encapsulates 043 the inter-dependency between agents but struggles to be scalable to an increasing number of agents, 044 which leads to the exponential surge in spatial complexity (Hernandez-Leal et al., 2020; Nguyen et al., 045 2020). Conversely, applying a *decentralized* world model to approximating the local dynamics of 046 each agent mitigates the scalability issue yet incurs non-stationarity, as unexpected interventions from 047 other agents may occur in each agent's individual environment (Oliehoek et al., 2016). Furthermore, 048 beyond these unique challenges inherent in modeling multi-agent dynamics, existing model-based MARL approaches (Willemsen et al., 2021; Egorov & Shpilman, 2022; Xu et al., 2022) excessively neglect the fact that the policy learned in imaginations of the world model heavily relies on the 051 quality of imagined trajectories (Micheli et al., 2023). It thereby necessitates accurate long-term prediction, especially with respect to the non-stationary local dynamics. Inspired by the capabil-052 ity of Transformer (Vaswani et al., 2017) in modeling complex discrete sequences and long-term dependency (Brown et al., 2020; Devlin et al., 2019; Micheli et al., 2023), we seek to construct a

Transformer-based world model within the multi-agent context for *decentralized* local dynamics together with *centralized* feature aggregation, combining the benefits of two distinctive designs.

In this paper, we introduce MARIE (Multi-Agent auto-Regressive Imagination for Efficient learning), the first Transformer-based multi-agent world model for sample-efficient policy learning. Specifically, the highlights of this paper are:

- 1. To tackle the inherent challenges within the multi-agent context, we build an effective world model via scalable *decentralized* dynamics modeling and essential *centralized* representation aggregating, which mirrors the principle of Centralized Training and Decentralized Execution.
- 2. To enable accurate and consistent long-term imaginations from the non-stationary local dynamics, we cast the *decentralized* dynamics learning as sequence modeling over discrete tokens by leveraging highly expressive Transformer architecture as the backbone. In particular, we successfully present the first Transformer-based world model for multi-agent systems.
 - 3. While it remains open for how to effectively enable *centralized* representation with the Transformer as the backbone, we achieve it by innovatively introducing a Perceiver Transformer (Jaegle et al., 2021) for efficient global information aggregation across all agents.
 - 4. Experiments on the Starcraft Multi-Agent Challenge (SMAC) benchmark in low data regime and additional experiments on MAMujoco show MARIE outperforms both model-free and existing model-based MARL methods w.r.t. both sample efficiency and overall performance and demonstrate the effectiveness of MARIE.
- 075 076 077

079

060

061

062

063

064

065

067

068

069

070

071

073

2 RELATED WORKS AND PRELIMINARIES

Multi-Agent Reinforcement Learning. In a model-free setting, a typical approach for cooperative MARL is centralized training and decentralized execution (CTDE), which tackles the scalability 081 and non-stationarity issues in MARL. During the training phase, it leverages global information 082 to facilitate agents' policy learning; while during the execution phase, it blinds itself and has only 083 access to the partial observation around each agent for multi-agent decision-making. Model-free 084 MARL methods with this paradigm can be divided into 2 categories: value-based (Sunehag et al., 085 2018; Rashid et al., 2018; Son et al., 2019; Wang et al., 2021) and policy-based (Lowe et al., 2017; Foerster et al., 2018; Iqbal & Sha, 2019; Ryu et al., 2020; Liu et al., 2020; Kuba et al., 2021; 087 Peng et al., 2021; Yu et al., 2022; Zhang et al., 2024b;a). In contrast to model-free approaches, 880 model-based MARL algorithms remain fairly understudied. MAMBPO (Willemsen et al., 2021) incorporates MBPO-style (Janner et al., 2019) techniques into multi-agent policy learning under the CTDE framework. Tesseract (Mahajan et al., 2021) introduces the tensorised Bellman equation and evaluates the Q-value function using Dynamic Programming (DP) together with an estimated 091 environment model. Similar to our setting where agents learn inside of an approximate world model, 092 MAMBA (Egorov & Shpilman, 2022) integrates the backbone proposed in DreamerV2 (Hafner et al., 2021) with an attention mechanism across agents to sustain an effective world model in environments 094 with an arbitrary number of agents, which leads to notably superior sample efficiency to existing 095 model-free approaches. In terms of model-based algorithm coupled with planning, MAZero (Liu 096 et al., 2024) expands the MCTS planning-based Muzero (Schrittwieser et al., 2020) framework to 097 the model-based MARL settings. However, learning-based or planning-based policies in these two 098 approaches are both overly coupled with their world models, downgrading their inference efficiency 099 and further limiting expansion in combinations with other popular model-free approaches. To the best of our knowledge, we are the first to expand the Transformer backbone-based world model within the 100 multi-agent context. 101

Learning behaviors within the imagination of world models. The Dyna architecture (Sutton, 1991)
 first emphasizes the utility of an estimated dynamics model in facilitating the training of the value
 function and policy. Inspired by the cognitive system of human beings, the concept of world model
 (Ha & Schmidhuber, 2018) is initially introduced by composing a variational Auto-Encoder (VAE)
 (Kingma & Welling, 2014) and a recurrent network to mimic the complete environmental dynamics,
 then an artificial agent is trained entirely inside the hallucinated imagination generated by the world
 model. SimPLe (Łukasz Kaiser et al., 2020) shows that a PPO policy (Schulman et al., 2017) learned



Figure 1: Overview of the proposed world model architecture in MARIE. VQ-VAE (left) maps 124 local observations o^i of each agent i into discrete latent codes $(x_1^i, ..., x_K^i)$, where (E, D, \mathcal{Z}) is 125 shared across all agents. Together with discrete actions, this process forms local discrete sequences 126 $(..., x_{t,1}^i, ..., x_{t,K}^i, a_t^i, ...)$ of each agent. Then the Perceiver (*right*) performs aggregation of joint dis-127 crete sequences of all agents $(x_{t,1}^1, ..., x_{t,K}^1, a_t^1, ..., x_{t,1}^n, ..., x_{t,K}^n, a_t^n)$ independently at each timestep t, and inserts the aggregated global representations $(e_t^1, e_t^2, ..., e_t^n)$ into original local discrete se-128 129 quences respectively. The resulting sequences $(..., x_{t,1}^i, ..., x_{t,K}^i, a_t^i, e_t^i...)$ contain rich information 130 between transitions in local dynamics and are fed into the shared Transformer (middle), which learns 131 observation token predictions in an autoregressive manner. Predictions of individual reward r_t^i and 132 discount γ_t^i at timestep t are computed based on all historical sequence $(x_{< t,1}^i, ..., x_{< t,K}^i, a_{< t}^i, e_{< t}^i)$. 133

in a predictive model deliverer a super-human performance in Atari domains. Dreamer (Hafner et al., 135 2020) builts the world model upon a Recurrent State Space Model (RSSM) (Hafner et al., 2019) that 136 combines the deterministic latent state with the stochastic latent state to allow the model to not only 137 capture multiple futures but also remember information over multi-steps. DreamerV2 (Hafner et al., 138 2021) further demonstrates the advantage of discrete latent states over Gaussian states. For MARL, 139 MAMBA (Egorov & Shpilman, 2022) extends DreamerV2 to multi-agent contexts by using RSSM, 140 underscoring the potential of multi-agent learning in the imagination of world models. Recently, 141 motivated by the success of the Transformer (Vaswani et al., 2017), TransDreamer (Chen et al., 2022) 142 and TWM (Robine et al., 2023) explored variants of DreamerV2, wherein the backbones of the world 143 model were substituted with Transformers. Instead of incorporating deterministic and stochastic 144 latent states, IRIS (Micheli et al., 2023) applies the Transformer to directly modeling sequences of 145 observation tokens and actions of single-agent RL and achieves impressive results on Atari-100k. In contrast, the proposed MARIE concentrates on establishing effective Transformer-based world 146 models in multi-agent contexts with shared dynamics and global representations. 147

148 Preliminaries. We focus on fully cooperative multi-agent systems where all agents share a team re-149 ward signal. We formulate the system as a decentralized partially observable Markov decision process 150 (Dec-POMDP) (Oliehoek et al., 2016), which can be described by a tuple $(\mathcal{N}, \mathcal{S}, \mathcal{A}, P, R, \Omega, \mathcal{O}, \gamma)$. $\mathcal{N} = \{1, ..., n\}$ denotes a set of agents, \mathcal{S} is the finite global state space, $\mathcal{A} = \prod_{i=1}^{n} \mathcal{A}^{i}$ is the product of finite actions spaces of all agents, i.e., the joint action space, $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ is the global 151 152 transition probability function, $R: S \times A \to \mathbb{R}$ is the shared reward function, $\Omega = \prod_{i=1}^{n} \Omega^{i}$ is the 153 product of finite observation spaces of all agents, i.e., the joint observation space, $\mathcal{O} = \{\mathcal{O}^i, i \in \mathcal{N}\}$ 154 is the set of observing functions of all agents. $\mathcal{O}^i: \mathcal{S} \to \tilde{\Omega}^i$ maps global states to the observations for 155 agent i, and γ is the discount factor. Given a global state s_t at timestep t, agent i is restricted to obtain-156 ing solely its local observation $o_t^i = \mathcal{O}^i(s_t)$, takes an action a_t^i drawn from its policy $\pi^i(\cdot | o_{\leq t}^i)$ based 157 on the history of its local observations $o_{\leq t}^i$, which together with other agents' actions gives a joint 158 action $a_t = (a_t^1, ..., a_t^n) \in \mathcal{A}$, equivalently drawn from a joint policy $\pi(\cdot | o_{\leq t}) = \prod_{i=1}^n \pi^i(\cdot | o_{< t}^i)$. 159 Then the agents receive a shared reward $r_t = R(s_t, a_t)$, and the environment moves to next state s_{t+1} with probability $P(s_{t+1}|s_t, a_t)$. The aim of all agents is to learn a joint policy π that maximizes 160 161 the expected discounted return $J(\boldsymbol{\pi}) = \mathbb{E}_{s_0, \boldsymbol{a}_0, \dots \sim \boldsymbol{\pi}} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, \boldsymbol{a}_t) \right].$

¹⁶² 3 METHODOLOGIES

163 164

Our approach comprises three typical parts: (i) collecting experience by executing the policy, (ii) learning the world model from the collected experience, and (iii) learning the policy via imagination inside the world model. Throughout the process, the historical experiences stored in the replay buffer are used for training the world model only, while policies are learned from unlimited imagined trajectories from the world model. In the following, we first describe three core components of our world model in §3.1 and §3.2, and give an overview of the proposed world model in Fig. 1. Then we describe the policy-learning process inside the world model in §3.3. The comprehensive details of the model architecture and hyperparameter are provided in §A.

 $\tau^{i} = (o_{1}^{i}, a_{1}^{i}, \dots, o_{t}^{i}, a_{t}^{i}, \dots, o_{T}^{i}, a_{T}^{i}).$

172 173

174

3.1 DISCRETIZING OBSERVATION

We consider a trajectory τ^i of agent *i* consists of *T* local observations and actions, as

175 176

195

196

201

177 To utilize the expressive Transformer architecture, we need to express the trajectory into a discrete 178 token sequence for modeling. Accounting for continuous observations, a prevalent but naive approach 179 for discretization involves discretizing the scalar into one of m fixed-width bins in each dimension 180 independently (Janner et al., 2021). However, when faced with a higher dimension of the observation, 181 such discretization would encode the observation with more tokens, leading to higher computational 182 complexity of the later sequence modeling via the Transformer, which necessitates an approach that 183 uses a discrete codebook of learned compact representations. To this end, we employ the idea from neural discrete representation learning (van den Oord et al., 2017), and learn a Vector Quantised-185 Variational AutoEncoder (VQ-VAE) to play a role that resembles the tokenizer in Natural Language Processing (Devlin et al., 2019; Brown et al., 2020). The VQ-VAE is composed of an encoder E, a decoder D, and a codebook Z. We define the discrete codebook $Z = \{z_j\}_{j=1}^N \subset \mathbb{R}^{n_z}$, where N is the size of the codebook and n_z is the dimension of codes. The encoder E takes an observation $o^i \in \mathbb{R}^{n_{obs}}$ as input and outputs a $K n_z$ -dimensional latents $\hat{z}^i \in \mathbb{R}^{K \times n_z}$ reshaped from the direct outputs of encoder. Subsequently, the takes $\{z_i\}_{j=1}^K$ 187 188 189 outputs of encoder. Subsequently, the tokens $\{x_k^i\}_{k=1}^K \in \{0, 1, ..., N-1\}^K$ for representing o^i is obtained by a nearest neighbour look-up using the codebook \mathcal{Z} where $x_k^i = \arg\min_j ||\hat{z}_k^i - z_j||$. 190 191 Then the decoder $D: \{0, 1, ..., N-1\}^K \to \mathbb{R}^{n_{obs}}$ converts K tokens back into an reconstructed 192 observation \hat{o}^i . By learning this discrete codebook, we compress the redundant information via a 193 succinct sequence of tokens, which helps improve sequence modeling. See §4.2 for a discussion. 194

3.2 MODELING LOCAL DYNAMICS WITH GLOBAL REPRESENTATIONS

Here, we consider discrete actions like those in SMAC, and the continuous actions can also be discretized by splitting the value in each dimension into fixed bins (Janner et al., 2021; Brohan et al., 2023). Therefore, a trajectory τ^i of agent *i* can be treated as a sequence of tokens,

$$\tau^{i} = (\dots, o_{t}^{i}, a_{t}^{i}, \dots) = (\dots, x_{t,1}^{i}, x_{t,2}^{i}, \dots, x_{t,K}^{i}, a_{t}^{i}, \dots)$$
(1)

where $x_{t,j}^i$ is the *j*-th token of the observation of agent *i* at timestep *t*. Given arbitrary sequences of observation and action tokens in Eq. (1), we try to learn over discrete multimodal tokens.

The world model consists of a tokenizer to discrete the local observation, a Transformer to learn the local dynamics, an agent-wise representation aggregation module, and predictors for the reward and discount. The Transformer ϕ predicts the future local observation $\{\hat{x}_{t+1,j}^i\}_{j=1}^K$, the future individual reward \hat{r}_t^i and discount $\hat{\gamma}_t^i$, based on the agent's individual historical observation-action history $(x_{<t}^i, a_{<t}^i)$ and aggregated global feature e_t^i of the agent. The modules are shown in Eqs. (2)–(5).

Transition:
$$\hat{x}_{t+1,.}^{i} \sim p_{\phi}(\hat{x}_{t+1,.}^{i} | x_{\leq t,.}^{i}, a_{\leq t}^{i}, e_{\leq t}^{i})$$
 with $\hat{x}_{t+1,k}^{i} \sim p_{\phi}(\hat{x}_{t+1,k}^{i} | x_{\leq t,.}^{i}, a_{\leq t}^{i}, e_{\leq t}^{i}, x_{t+1,
(2)$

213 Reward:
$$\hat{r}_{t}^{i} \sim p_{\phi}(\hat{r}_{t}^{i}|x_{\leq t}^{i}, a_{\leq t}^{i}, e_{\leq t}^{i})$$
 (3)
214 D: $\hat{r}_{t}^{i} \sim p_{\phi}(\hat{r}_{t}^{i}|x_{\leq t}^{i}, a_{\leq t}^{i}, e_{\leq t}^{i})$ (3)

215 Discount:
$$\gamma_t^* \sim p_\phi(\gamma_t^* | x_{\le t,.}^*, a_{\le t}^*, e_{\le t}^*)$$
 (4)

Aggregation:
$$(e_t^1, e_t^2, ..., e_t^n) = f_\theta(x_{t,1}^1, x_{t,2}^1, ..., x_{t,K}^1, a_t^1, ..., x_{t,1}^n, x_{t,2}^n, ..., x_{t,K}^n, a_t^n)$$
 (5)

216 **Transition Prediction.** In the transition prediction in Eq. (2), the k-th observation token is additionally 217 conditioned on the tokens that were already predicted $x_{t+1,<k}^i \triangleq (x_{t+1,1}^i, x_{t+1,2}^i, ..., x_{t+1,k-1}^i),$ 218 ensuring the autoregressive token prediction to facilitate modeling over the trajectory sequence. 219 Inter-step auto regression is as intuitive as predicting the future based on all information in the past 220 while intra-step auto regression can be interpreted as learning how to compose the language provided 221 by VQ-VAE to correctly express the observation within a certain timestep, since the tokens for 222 encoding observations can be viewed as a special inner language like the human's.

223 **Discount and Reward Prediction.** The discount predictor outputs a Bernoulli likelihood and lets us 224 estimate the probability of an individual agent's episode ending when learning behaviors from model 225 predictions. And we simply adopt a smooth L1 loss for training the prediction of reward. 226

Agent-wise Aggregation. Due to the partial environment, the non-stationarity issue stems from 227 the sophisticated agent-wise inter-dependency on local observations generation. To address it, we 228 introduce a Perceiver (Jaegle et al., 2021) to perform agent-wise representation aggregation which 229 plays a similar role to communication. To sustain the decentralized manner in transition prediction, 230 we hope every agent can possess its own inner perception of the whole situation. Nonetheless, with 231 discrete representation for local observation, the observation-action pair of agent i at timestep t is 232 projected into a sequence $(x_{t,1}^i, x_{t,2}^i, ..., x_{t,K}^i, a_t^i)$ of length K+1. It leads to a joint observation-action 233 sequence of length n(K+1) at a timestep, which linearly scales with the number of agents. 234

A naive approach for extracting aggregated feature for each agent is using self-attention (Egorov & 235 Shpilman, 2022; Liu et al., 2024) which takes as input this sequence of length n(K+1) and outputs 236 a sequence of the same length containing aggregated features of all agents, described as 237

$$(x_{t,1}^1, ..., x_{t,K}^1, a_t^1, ..., x_{t,1}^n, ..., x_{t,K}^n, a_t^n) \xrightarrow{\text{Self-Attention}}_{\text{Aggregating}} (e_{t,1}^1, ..., e_{t,K}^1, e_{t,K+1}^1, ..., e_{t,1}^n, ..., e_{t,K}^n, e_{t,K+1}^n).$$

240 where $e_{t,j}^i$ is the j-th aggregated feature for agent i at timestep t. However, when composing the 241 informative sequence of local trajectories by insert these aggregated features into the sequence of 242 length H(K + 1) in Eq. (1), the length of local sequence involving aggregated features would be 243 twice longer, i.e., 2H(K+1). Due to the quadratic computational complexity of Transformer, it may 244 hinder the efficient sequence modeling over this sequence. 245

To this end, we choose the Perceiver as the agent-wise representation aggregation module, which 246 excels at dealing with the case that the size of inputs scales linearly and then generates a compact 247 output sequence. Equipped with a flexible querying mechanism and self-attention mechanism, the Per-248 ceiver aggregates the joint representation sequence $(x_{t,1}^1, x_{t,2}^1, ..., x_{t,K}^1, a_t^1, ..., x_{t,1}^n, x_{t,2}^n, ..., x_{t,K}^n, a_t^n)$ 249 of length n(K+1) into a sequence of n features $(e_t^1, e_t^2, ..., e_t^n)$, 250

251 252

253

254

255

256

259

261 262

238 239

$$(x_{t,1}^1,...,x_{t,K}^1,a_t^1,...,x_{t,1}^n,...,x_{t,K}^n,a_t^n) \xrightarrow{\operatorname{Perceiver}}_{\operatorname{Aggregating}} (e_t^1,e_t^2,...,e_t^n)$$

where each feature e_t^i serves as an intrinsic global abstraction of the environmental contexts perceived from agent *i*'s viewpoint. By introducing Perceiver, we provide a feasible solution for reducing the modeling complexity when using transformer-based local dynamics.

Overall Learning Objective. The world model ϕ is trained with trajectory segments of a fixed 257 horizon H sampled from the replay buffer \mathcal{D} in a self-supervised manner. The transition predictor, 258 discount predictor, and reward predictor are optimized to maximize the log-likelihood of their corresponding targets: 260

$$\mathcal{L}_{\text{Dyn}}(\phi,\theta) = \mathbb{E}_{i\sim\mathcal{N}} \mathbb{E}_{\tau^{i}\sim\mathcal{D}} \left[\sum_{t=1}^{H} -\underbrace{\log p_{\phi}(r_{t}^{i}|x_{\leq t,\cdot}^{i},a_{\leq t}^{i},e_{t}^{i})}_{\text{reward loss}} -\underbrace{\log p_{\phi}(\gamma_{t}^{i}|x_{\leq t,\cdot}^{i},a_{\leq t}^{i},e_{t}^{i})}_{\text{discount loss}} -\underbrace{\left(\sum_{k=1}^{K} \log p_{\phi}(x_{t+1,k}^{i}|x_{\leq t,\cdot}^{i},a_{\leq t}^{i},e_{t}^{i},x_{t+1,(6)$$

where
$$(e_t^1, e_t^2, ..., e_t^n) = f_{\theta}(x_{t,1}^1, x_{t,2}^1, ..., x_{t,K}^1, a_t^1, ..., x_{t,1}^n, x_{t,2}^n, ..., x_{t,K}^n, a_t^n), \forall t.$$

We jointly minimize this loss function in Eq. (6) with respect to the model parameters of local 269 dynamics (i.e., ϕ) and global representation (i.e., θ) using the Adam optimizer (Kingma & Ba, 2015).



Figure 2: Imagination procedure in MARIE. We unroll the imagination of all agents $\{1, ..., n\}$ in parallel. Initially, each agent's observation is derived from a joint observation sampled from a replay buffer. A policy, depicted in red arrows, generates actions based on reconstructed observations. Then, the Perceiver integrates joint actions and observations into global representations from each agent, appending them to each agent's local sequence. The Transformer then predicts individual rewards and discounts, depicted by green and purple arrows respectively, while generating next observation tokens for each agent in an autoregressive manner, shown by blue arrows. This parallel imagination iterates for H steps. The policies $\pi_{\psi}^{1:n}$ are exclusively trained using imagined trajectories.

3.3 LEARNING BEHAVIOURS IN IMAGINATION

290 We utilize the Actor-Critic framework to learn the behavior of each agent, where the actor and critic are parameterized by ψ and ξ , respectively. In the following, we take agent i as an exemplar case 291 for clarity and omit the superscript for denoting the index of the agent to avoid potential confusion. 292 Benefited from the shared local dynamics, the local trajectories of all agents are imagined in parallel, 293 as illustrated in Fig. 2. At timestep t, the actor takes a reconstructed observation \hat{o}_t as input, and samples an action $a_t \sim \pi_{\psi}(a_t | \hat{o}_t)$. The world model then predicts the individual reward \hat{r}_t , individual 295 discount $\hat{\gamma}_t$ and next local observation $\hat{\sigma}_{t+1}$. Starting from initial observations sampled from the 296 replay buffer, this imagination procedure is rolled out for H steps. To stimulate long-horizon behavior 297 learning, the critic accounts for rewards beyond the fixed imagination horizon and estimates the individual expected return $V_{\xi}(\hat{o}_t) \simeq \mathbb{E}_{\pi_{\psi}}[\sum_{l \ge t} \gamma^{l-t} \hat{r}_l].$ 298 299

In our approach, we train the actor and critic in a MAPPO-like (Yu et al., 2022) manner. Unlike other 300 CTDE model-free approaches that require a global oracle state from the environment, we cannot 301 obtain the oracle state from the world model, and only the predicted observations of each agent are 302 available. To approximate the oracle information in critic training, we enhance each agent's critic 303 with the capability to access the observations of other agents. Since the actor and critic only rely on 304 the reconstructed observations, decoupling from the inner hidden states of the Transformer-based 305 world model, we allow fast inference in the environment without the participation of the world model. 306 It is important for the deployment of policies learned with data-efficient imagination in real-world 307 applications. λ -target in Dreamer (Hafner et al., 2020) is used to updated the value function. The 308 details of behavior learning objectives and algorithmic description of MARIE are presented in §B and §I, respectively. 309

310 311

4 EXPERIMENTS

312 313

We consider the most common benchmark – StarCraftII Multi-Agent Challenge (SMAC) (Samvelyan et al., 2019) for evaluating our method. To highlight the sample efficiency brought by model-based imagination, we adopt a low data regime that resembles a similar setting in single-agent Atari domain (Łukasz Kaiser et al., 2020). Additional experiment results on MAMujoco (Peng et al., 2021) (i.e., continuous action space case) is provided in §E.1.

318 319

320

4.1 EXPERIMENT SETUP AND EVALUATIONS

StarCraftII Multi-Agent Challenge. SMAC (Samvelyan et al., 2019), a suite of cooperative multi agent environments based on StarCraft II, consists of a set of StarCraft II scenarios. Each scenario
 depicts a confrontation between two armies of units, one of which is controlled by the built-in game
 AI and the other by our algorithm. The initial position, number, and type of units in each army

339

340



Figure 3: Curves of evaluation win rate for methods in 8 chosen SMAC maps. See Table 1 for win rates. Y axis: win rate; X axis: number of steps taken in the real environment. MARIE demonstrates superior performance and sample efficiency across almost all scenarios.

341 varies from scenario to scenario, as does the presence or absence of elevated or impassable terrain. And the goal is to win the game within the pre-specified time limit. SMAC emphasizes mastering 342 micromanagement techniques across multiple agents to achieve effective coordination and overcome 343 adversaries. This necessitates both sufficient exploration and appropriate credit assignment for each 344 agent's action. Another notable property of SMAC is that not all actions are accessible during 345 decision-making of each agent, which requires world models to possess an in-depth comprehension 346 of the underlying game mechanics so as to consistently provide valid available action mask estimation 347 within the imagination horizon. Thus, in this benchmark, we additionally add one more head for the 348 prediction of available action mask. During the imagination of MARIE, the available action mask 349 is estimated by this head, instead of being generated manually according to the meaning of each 350 element in the reconstructed observation. The latter introduces too much prior knowledge about 351 StarCraft and can be considered as benchmark hacking.

352 **Experimental Setup.** We choose 13 representative scenarios from SMAC that includes three levels 353 of difficulty – *Easy*, *Hard*, and *SuperHard*. Specific chosen scenarios can be found in Table 1. In 354 terms of different levels of difficulty, we adopt a similar setting akin to that in (Egorov & Shpilman, 355 2022) and restrict the number of samples from the real environment to 100k for Easy scenarios, 200k 356 for *Hard* scenarios and 400k for *SuperHard* scenarios, to establish a low data regime in SMAC. 357 We compare MARIE with three strong model-free baselines – MAPPO (Yu et al., 2022), QMIX 358 (Rashid et al., 2018) and OPLEX (Wang et al., 2021), and two strong model-based baselines with the same policy learning paradigm as ours – MBVD (Xu et al., 2022) and MAMBA (Egorov & 359 Shpilman, 2022) on SMAC benchmark. Specially, as a multi-agent variant of DreamerV2 (Hafner 360 et al., 2021), MAMBA achieves powerful sample efficiency in various SMAC scenarios via learning 361 in imagination. For each random seed, we compute the win rate across 10 evaluation games at fixed 362 intervals of environmental steps. The hyperparameters of MARIE and other baselines are listed in \$D and \$H. Particularly, the hyperparameters of model-free baselines in low data regime are directly 364 referred to Egorov & Shpilman (2022) and Liu et al. (2024).

Main Results. Overall, we find MARIE achieves significantly better sample efficiency and a higher 366 win rate compared with other strong baselines. We report the averaged win rates over four seeds 367 in Table 1 and provide additional learning curves of several chosen scenarios, shown as Fig. 3. As 368 presented in Table 1 and Fig. 3, MARIE demonstrates superior performance and sample efficiency 369 across almost all scenarios. The improvements in sample efficiency and performance become 370 particularly pronounced with increasing difficulty of scenarios, especially compared to MAMBA that 371 adopts RSSM as the backbone for the world model. We attribute such results to the model capability 372 of the Transformer in local dynamics modeling and global feature aggregation. Benefiting from more 373 powerful strength in modeling sequences, the Transformer-based world model can generate more 374 accurate and consistent imaginations than those relying on the recurrent backbone, which facilitates 375 better policy learning within the imagination of the world model. While the scenarios become harder, e.g. $3s_vs_5z$, our world model can address the challenge of learning more intricate underlying 376 dynamics and further large quantities of accurate imaginations, thereby significantly outperforming 377 other baselines on these scenarios. Moreover, a special scenario $2c_{-}v_{s}_{-}64zg$ deserves attention, which

396

397

398

399

400

401

402

403

404

405

406 407

Table 1: Mean evaluation win rate and standard deviation on 13 SMAC maps for different methods over 4 random seeds. We bold the values of the maximum and highlight them with blue color.

	-	~			Me	thods		
Maps	Difficulty	Steps	MARIE (Ours)	MAMBA (Egorov & Shpilman, 2022)	MAPPO (Yu et al., 2022)	QMIX (Rashid et al., 2018)	QPLEX (Wang et al., 2021)	MBVD (Xu et al., 2022)
1c3s5z			85.0 (9.4)	77.7(15.3)	18.4(11.0)	43.6(29.2)	68.3(7.4)	60.9(11.4)
2m_vs_1z			95.5 (7.9)	95.5(2.3)	86.7(3.2)	70.3(14.8)	84.8(10.8)	36.7(24.5)
2s_vs_1sc			96.9(7.1)	95.0(7.1)	100.0(0.0)	0.0(0.0)	15.7(19.5)	8.7(14.8)
2s3z		10017	80.5(9.3)	71.6(12.7)	31.2(12.9)	37.7(15.5)	50.2(8.4)	53.4(4.1)
3m	Easy	100K	99.5 (0.4)	87.7(7.1)	80.5(12.8)	54.4(22.7)	88.7(6.9)	73.9(6.9)
3s_vs_3z			98.9 (1.5)	89.3(10.1)	1.2(1.3)	0.0(0.0)	0.0(0.0)	0.0(0.0)
3s_vs_4z			73.0(6.2)	29.3(12.3)	0.0(0.0)	0.0(0.0)	0.0(0.0)	0.0(0.0)
8m			88.0 (3.9)	65.0(7.7)	70.3(19.5)	69.5(12.8)	83.4(6.4)	74.7(9.7)
MMM			87.6(3.0)	50.2(27.6)	5.5(4.5)	31.1(17.3)	69.3(35.1)	20.5(2.1)
so_many_baneling			94.8 (5.9)	91.6(4.1)	43.8(15.0)	20.0(8.9)	32.2(6.1)	15.0(10.4)
3s_vs_5z	11	2001/	78.4(11.2)	13.4(14.0)	0.0(0.0)	0.0(0.0)	0.0(0.0)	0.0(0.0)
2c_vs_64zg	nara	200 K	25.9 (14.3)	9.8(8.7)	7.8(10.2)	0.5(0.5)	0.1(0.1)	0.2(0.4)
corridor	SuperHard	400K	71.0(13.8)	26.5(15.2)	0.4(0.7)	0.0(0.0)	0.0(0.0)	0.0(0.0)



Figure 4: Ablation on what manner to integrate into the design of the world model. Decentralized Manner denotes the standard implementation of MARIE, while Centralized Manner denotes that the world model is designed for learning the joint dynamics of all agents over the joint trajectory. Sample efficiency of the *centralized* variant encounters a significant drop due to the scalability issue while MARIE is robust to scenarios with various number of agents.

features only 2 agents but with a considerably large action space of up to 70 discrete actions for each 408 agent. Although the performance of MARIE in $2c_{vs}$ -64zg suffers a relative large variance due to the 409 overly large action space, MARIE achieves a remarkably non-trivial mean win rate just via learning 410 in the imagination. Note that it is easy for the world model to generate ridiculous estimated available 411 action masks without understanding the mechanics behind this scenario, further leading to invalid 412 or even erroneous policy learning in the imaginations of the world model. The performance gap on 413 $2c_{vs}$. 64zg proves that our Transformer-based world model has higher prediction accuracy and a 414 deeper understanding of the underlying mechanics.

415 416

418

4.2 ABLATION STUDIES 417

Incorporating CTDE principle with the design of the world model makes MARIE scalable 419 and robust to different number of agents. We compare our method with a centralized variant 420 of our method, wherein the world model learns the joint dynamics of all agents together over the 421 joint trajectory $\tau = (\dots, o_t^1, o_t^2, \dots, o_t^n, a_t^1, a_t^2, \dots, a_t^n, \dots)$. Given that τ already contains the joint 422 observations and actions, we disable the aggregation module in this *centralized* variant. As illustrated 423 in Figure 4, our comparisons span scenarios involving 2 to 7 agents. When the number of agents is 424 small enough, reducing the multi-agent system to a single-agent one over the joint observation and 425 action space would not cause a prominent scalability issue, as indicated by the result in $2s_v s_s 1s_c$. 426 However, the scalability issue is exacerbated by a growing number of agents. In scenarios featuring 427 more than 3 agents, the sample efficiency of the *centralized* variant encounters a significant drop, 428 suffering from the exponential surge in spatial complexity of the joint observation-action space. Furthermore, with equal prediction horizons, the parameter amounts in the *centralized* variant is 429 increased by a factor of 4 or larger. And to achieve the same number of environment steps, the 430 centralized variant demands over twice the original computational time. Instead, with decentralized 431 local dynamics and aggregated global features, MARIE delivers stable and superior sample efficiency.

433

434

435

437

438

439

440

441

446

447

448

449

450

451

452

453

454

455

456

457

458 459 460



Figure 5: Comparisons between MARIE with and without the usage of the aggregation module. Local dynamics struggles to infer accurate future local observations without agent-wise aggregation.



Figure 6: Ablation on the type of discretization for local observations. Tokenizer denotes the standard implementation of MARIE; Bins Discretization denotes the variant of MARIE where the n_{obs} -dimensional observation discretization is performed by projecting the value into one of m fixed-width bins in each dimension independently. X-axis: cumulative run time of algorithms in the same platform. VQ-VAE encapsulates local observations within a succinct sequence of tokens, computationally efficiently promoting the learning of the Transformer-based world model.

461 Agent-wise aggregation helps MARIE capture the sophisticated inter-dependency on the genera-462 tion of each agent's local observation. To study the influence of agent-wise aggregation, we conduct 463 ablation experiments on the aggregation module over scenarios where the number of agents gradually 464 increases. As shown in Fig. 5, in the 3-agents scenario (e.g., $3s_vs_3z$), the correlation among each 465 agent's local observation tends to be negligible. Therefore, the nearly independent generation of each agent's local observation without any aggregated global feature still leads to performance comparable 466 to that of standard implementation. But as more agents get involved, the inter-dependency becomes 467 dominant. Lacking the global features derived from agent-wise aggregation, the shared Transformer 468 struggles to infer accurate future local observations, thus hindering policy learning in the imaginations 469 of the world model and resulting in notable degradation in the win rate evaluation. 470

471 VO-VAE encapsulates local observations within a succinct sequence of tokens, promoting the learning of the Transformer-based world model and effectively improving algorithm perfor-472 **mance.** Compared to VQ-VAE that discretizes each observation to K tokens from \mathcal{Z} , perhaps a more 473 naive tokenizer is projecting the value in each dimension into one of m fixed-width bins (Janner 474 et al., 2021), resulting in a n_{obs} -long token sequence for each observation, which we term *Bins* 475 *Discretization.* We set the number of bins m equal to the size of codebook $|\mathcal{Z}|$ and compare these two 476 types of tokenizers in different environments with various n_{obs} . As shown in Fig. 6, the performance 477 of the two tokenizers are comparable only in $2s_vs_lsc$ where n_{obs} is close to 16. Even worse, *Bins* 478 Discretization experiences a pronounced decline as $n_{\rm obs}$ increases in more complex environments 479 (e.g., $3s_v v_s 4z$) under identical training durations. We hypothesize that for a single local observation, 480 a n_{obs}-token-long verbose sequence yielded by *Bins Discretization* contains more redundant infor-481 mation compared to VQ-VAE that learns a more compact tokenizer through reconstruction This not 482 only renders the token sequences of *Bins Discretization* obscure and challenging to comprehend, but 483 also results in an increase in model parameter amounts, being more computationally costly. Due to these two factors, *Bins Discretization* exhibits a notably slow convergence. Meanwhile, the result 484 in 2m_vs_1z indicates Bins Discretization may ignore the correlation of different dimensions, which 485 would be helpful in sequence modeling.



Figure 7: **Compounding model errors.** We compare the imagination accuracy of MARIE to that of MAMBA over the course of a planning horizon in $3s_{vs_{-}5z}$ scenario. MARIE has remarkably better error compounding with respect to prediction horizon than MAMBA.

4.3 MODEL ANALYSIS

Error Accumulation. A quantitative evaluation of the model's accumulated error versus prediction 501 horizon is provided in Fig. 7. Since learning the world model is tied to a progressively improving 502 policy both in MARIE and MAMBA, we separately use their final policies to sample 10 episodes for 503 fairness. We then compute L_1 errors per observation dimension between 1000 trajectory segments 504 randomly sampled from these 20 episodes and their imagined counterpart. The result in Fig. 7 505 suggests architecture differences play a large role in the world model's long-horizon accuracy. 506 This also provides additional evidence that policy learning can benefit from accurate long-term 507 imaginations, explaining MARIE's notable performance in the 3s_vs_5z scenario. More precisely, 508 lower generalization error between the estimated dynamics and true dynamics brings a tighter bound 509 between optimal policies derived from these two dynamics according to theoretical results (Janner 510 et al., 2019).

511 Attention Patterns. During model prediction, we delve into the attention maps inside the shared 512 Transformer and the cross attention maps in the Perceiver. Interestingly, we observe two distinct 513 attention patterns involved in the local dynamics prediction. One exhibits a Markovian pattern 514 wherein the observation prediction lays its focus mostly on the previous transition, while the other is 515 regularly striated wherein the model attends to specific tokens in multiple prior transitions. During 516 the agent-wise aggregation, we also identify two distinct patterns – *individuality* and *commonality* 517 among agents. Such diverse patterns in the Transformer and Perceiver may be pivotal for achieving 518 accurate and consistent imaginations of the sophisticated local dynamics. We refer to §C for further details and visualization results. 519

520 521

522 523

495

496

497 498 499

500

5 CONCLUSION AND LIMITATION

524 We have introduced a model-based multi-agent algorithm – MARIE, which utilizes a shared Trans-525 former as local dynamic model and a Perceiver as a global agent-wise aggregation module to construct a world model within the multi-agent context. By providing long-term imaginations with policy 526 learning, it significantly boosts the sample efficiency and improves final performance compared to 527 state-of-the-art model-free methods and existing model-based methods with same learning paradigm, 528 in the low data regime. But it should be also noticed that there are potential limitations on the 529 current evaluation on the main experiment with 4 limited seeds, e.g., the limitations of mean and 530 median scores (Agarwal et al., 2021). Thus, we also provide a standardized performance evaluation 531 following the protocol provided by Agarwal et al. (2021) in §G. To further deliver a rigorous statis-532 tical validation, evaluation with more seeds is definitely necessary. As the first Transformer-based 533 multi-agent world model for sample-efficient policy learning, we open a new avenue for combining 534 the powerful strength of the Transformer with sample-efficient MARL. Considering the notorious sample inefficiency in multi-agent scenarios, it holds important promise for application in many realistic multi-robot systems, wherein collecting tremendous samples for optimal policy learning is 537 costly and impractical due to the safety. While it has the great potential to bright the future towards achieving smarter multi-agent systems, there still exist limitations in MARIE. For instance, it would 538 suffer from much slower inference speed when used with a very long prediction horizons, due to the auto-regressive property.

540 REPRODUCIBILITY STATEMENT

For the implementation details, we provide the detailed instruction in §A. For the practical part, we give experiment setup in §4. The hyper-parameters and implementation details are given in §H. The code will be released publicly after the review process.

546 547 REFERENCES

548

549

550

580

581

582

583

592

- Rishabh Agarwal, Max Schwarzer, Pablo Samuel Castro, Aaron C Courville, and Marc Bellemare. Deep reinforcement learning at the edge of the statistical precipice. *Advances in neural information processing systems*, 34:29304–29320, 2021.
- Shaked Brody, Uri Alon, and Eran Yahav. How attentive are graph attention networks? *arXiv preprint arXiv:2105.14491*, 2021.

554 Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, 555 Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, Julian Ibarz, Brian Ichter, 556 Alex Irpan, Tomas Jackson, Sally Jesmonth, Nikhil Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Kuang-Huei Lee, Sergey Levine, Yao Lu, Utsav Malla, Deeksha 558 Manjunath, Igor Mordatch, Ofir Nachum, Carolina Parada, Jodilyn Peralta, Emily Perez, Karl 559 Pertsch, Jornell Quiambao, Kanishka Rao, Michael Ryoo, Grecia Salazar, Pannag Sanketi, Kevin Sayed, Jaspiar Singh, Sumedh Sontakke, Austin Stone, Clayton Tan, Huong Tran, Vincent Vanhoucke, Steve Vega, Quan Vuong, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna 561 Zitkovich. Rt-1: Robotics transformer for real-world control at scale. In Robotics: Science and 562 Systems (RSS), 2023. 563

564 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, 565 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel 566 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, 567 Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, 568 Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In Advances in Neural Information 569 Processing Systems, 2020. URL https://proceedings.neurips.cc/paper_files/ 570 paper/2020/file/1457c0d6bfcb4967418bfb8ac142f64a-Paper.pdf. 571

- 572
 573
 574
 Chang Chen, Yi-Fu Wu, Jaesik Yoon, and Sungjin Ahn. Transdreamer: Reinforcement learning with transformer world models. *arXiv preprint arXiv:2202.09481*, 2022.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep
 bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT*, 2019. doi: 10.18653/V1/N19-1423. URL https://doi.org/10.
 18653/v1/n19-1423.
 - Vladimir Egorov and Alexei Shpilman. Scalable multi-agent model-based reinforcement learning. In *Proceedings of the 21st International Conference on Autonomous Agents and Multiagent Systems*, 2022.
- Benjamin Ellis, Jonathan Cook, Skander Moalla, Mikayel Samvelyan, Mingfei Sun, Anuj Mahajan, Jakob Nicolaus Foerster, and Shimon Whiteson. SMACv2: An improved benchmark for cooperative multi-agent reinforcement learning. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023. URL https://openreview.net/forum?id=50jLGiJW3u.
- Jakob N. Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson.
 Counterfactual multi-agent policy gradients. In *Proceedings of the Thirty-Second AAAI Conference* on Artificial Intelligence, 2018.
- ⁵⁹³ David Ha and Jürgen Schmidhuber. Recurrent world models facilitate policy evolution. In *Advances in Neural Information Processing Systems*, 2018.

594 595 596 597	Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and James Davidson. Learning latent dynamics for planning from pixels. In <i>Proceedings of the 36th</i> <i>International Conference on Machine Learning</i> , Proceedings of Machine Learning Research. PMLR, 2019.
598 599	Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning
600 601	URL https://openreview.net/forum?id=S110TC4tDS.
602 603 604	Danijar Hafner, Timothy P Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering atari with discrete world models. In <i>International Conference on Learning Representations</i> , 2021. URL https://openreview.net/forum?id=0oabwyZbOu.
605 606 607	Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains through world models. <i>arXiv preprint arXiv:2301.04104</i> , 2023.
608 609 610 611	Nicklas Hansen, Hao Su, and Xiaolong Wang. Temporal difference learning for model predictive control. In <i>Proceedings of the 39th International Conference on Machine Learning</i> , Proceedings of Machine Learning Research. PMLR, 2022. URL https://proceedings.mlr.press/v162/hansen22a.html.
612 613 614 615	Nicklas Hansen, Hao Su, and Xiaolong Wang. TD-MPC2: Scalable, robust world models for continuous control. In <i>The Twelfth International Conference on Learning Representations</i> , 2024. URL https://openreview.net/forum?id=Oxh5CstDJU.
616 617	Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). arXiv preprint arXiv:1606.08415, 2016.
618 619 620 621	Pablo Hernandez-Leal, Bilal Kartal, and Matthew E. Taylor. A very condensed survey and critique of multiagent deep reinforcement learning. In <i>Proceedings of the 19th International Conference on Autonomous Agents and MultiAgent Systems</i> , 2020.
622 623 624	Shariq Iqbal and Fei Sha. Actor-attention-critic for multi-agent reinforcement learning. In <i>Proceedings</i> of the 36th International Conference on Machine Learning, Proceedings of Machine Learning Research. PMLR, 2019. URL http://proceedings.mlr.press/v97/iqbal19a.html.
625 626 627 628	Andrew Jaegle, Felix Gimeno, Andy Brock, Oriol Vinyals, Andrew Zisserman, and Joao Carreira. Perceiver: General perception with iterative attention. In <i>International conference on machine learning</i> , 2021.
629 630	Michael Janner, Justin Fu, Marvin Zhang, and Sergey Levine. When to trust your model: Model-based policy optimization. In <i>Advances in Neural Information Processing Systems</i> , 2019.
631 632 633	Michael Janner, Qiyang Li, and Sergey Levine. Offline reinforcement learning as one big sequence modeling problem. In <i>Advances in Neural Information Processing Systems</i> , 2021.
634 635	Andrej Karpathy. mingpt: A minimal pytorch re-implementation of the openai gpt (generative pretrained transformer) training. https://github.com/karpathy/minGPT, 2020.
636 637 638 639	Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio and Yann LeCun (eds.), <i>3rd International Conference on Learning Representations</i> , 2015. URL http://arxiv.org/abs/1412.6980.
640 641	Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. In 2nd International Conference on Learning Representations, 2014.
642 643 644 645	Jakub Grudzien Kuba, Ruiqing Chen, Muning Wen, Ying Wen, Fanglei Sun, Jun Wang, and Yaodong Yang. Trust region policy optimisation in multi-agent reinforcement learning. In <i>International Conference on Learning Representations</i> , 2021.
646 647	Qihan Liu, Jianing Ye, Xiaoteng Ma, Jun Yang, Bin Liang, and Chongjie Zhang. Efficient multi- agent reinforcement learning by planning. In <i>The Twelfth International Conference on Learning</i> <i>Representations</i> , 2024. URL https://openreview.net/forum?id=CpnKq3UJwp.

648 649 650 651	Yong Liu, Weixun Wang, Yujing Hu, Jianye Hao, Xingguo Chen, and Yang Gao. Multi-agent game abstraction via graph attention neural network. In <i>The Thirty-Fourth AAAI Conference on Artificial Intelligence</i> , 2020. doi: 10.1609/AAAI.V34I05.6211. URL https://doi.org/10.1609/aaai.v34i05.6211.
652 653 654 655	Ryan Lowe, Yi Wu, Aviv Tamar, Jean Harb, Pieter Abbeel, and Igor Mordatch. Multi-agent actor- critic for mixed cooperative-competitive environments. In <i>Proceedings of the 31st International</i> <i>Conference on Neural Information Processing Systems</i> , 2017.
656 657 658 659 660	Anuj Mahajan, Mikayel Samvelyan, Lei Mao, Viktor Makoviychuk, Animesh Garg, Jean Kossaifi, Shimon Whiteson, Yuke Zhu, and Animashree Anandkumar. Tesseract: Tensorised actors for multi-agent reinforcement learning. In <i>Proceedings of the 38th International Conference on</i> <i>Machine Learning</i> , Proceedings of Machine Learning Research. PMLR, 2021. URL https: //proceedings.mlr.press/v139/mahajan21a.html.
661 662 663	Vincent Micheli, Eloi Alonso, and François Fleuret. Transformers are sample-efficient world models. In <i>The Eleventh International Conference on Learning Representations</i> , 2023. URL https://openreview.net/forum?id=vhFulAcb0xb.
664 665 666	Thanh Thi Nguyen, Ngoc Duy Nguyen, and Saeid Nahavandi. Deep reinforcement learning for multiagent systems: A review of challenges, solutions, and applications. <i>IEEE Transactions on Cybernetics</i> , 50:3826–3839, 2020. doi: 10.1109/TCYB.2020.2977374.
668 669	Frans A Oliehoek, Christopher Amato, et al. A concise introduction to decentralized POMDPs, volume 1. <i>Springer</i> , 2016.
670 671 672 673	Bei Peng, Tabish Rashid, Christian Schroeder de Witt, Pierre-Alexandre Kamienny, Philip Torr, Wendelin Boehmer, and Shimon Whiteson. FACMAC: Factored multi-agent centralised policy gradients. In <i>Advances in Neural Information Processing Systems</i> , 2021. URL https:// openreview.net/forum?id=WxH774N0mEu.
674 675 676	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 2019.
677 678 679 680 681	Tabish Rashid, Mikayel Samvelyan, Christian Schroeder, Gregory Farquhar, Jakob Foerster, and Shi- mon Whiteson. QMIX: Monotonic value function factorisation for deep multi-agent reinforcement learning. In <i>Proceedings of the 35th International Conference on Machine Learning</i> , Proceedings of Machine Learning Research. PMLR, 2018. URL https://proceedings.mlr.press/ v80/rashid18a.html.
682 683 684	Jan Robine, Marc Höftmann, Tobias Uelwer, and Stefan Harmeling. Transformer-based world models are happy with 100k interactions. In <i>The Eleventh International Conference on Learning Representations</i> , 2023. URL https://openreview.net/forum?id=TdBaDGCpjly.
685 686 687	Heechang Ryu, Hayong Shin, and Jinkyoo Park. Multi-agent actor-critic with hierarchical graph attention network. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 2020.
688 689 690 691	Mikayel Samvelyan, Tabish Rashid, Christian Schroeder de Witt, Gregory Farquhar, Nantas Nardelli, Tim G. J. Rudner, Chia-Man Hung, Philip H. S. Torr, Jakob Foerster, and Shimon Whiteson. The starcraft multi-agent challenge. In <i>Proceedings of the 18th International Conference on</i> <i>Autonomous Agents and MultiAgent Systems</i> , 2019.
692 693 694 695	Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, L. Sifre, Simon Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, Timothy P. Lillicrap, and David Silver. Mastering atari, go, chess and shogi by planning with a learned model. <i>Nature</i> , 588:604 – 609, 2020.
696 697 698	John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. <i>arXiv preprint arXiv:1707.06347</i> , 2017.
699 700 701	Kyunghwan Son, Daewoo Kim, Wan Ju Kang, David Earl Hostallero, and Yung Yi. QTRAN: Learning to factorize with transformation for cooperative multi-agent reinforcement learning. In <i>Proceedings of the 36th International Conference on Machine Learning</i> , Proceedings of Machine Learning Research. PMLR, 2019. URL https://proceedings.mlr.press/v97/son19a.html.

702 703 704 705 706	Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi, Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z. Leibo, Karl Tuyls, and Thore Graepel. Value-decomposition networks for cooperative multi-agent learning based on team reward. In <i>Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems</i> , 2018.
707 708 709	Richard S Sutton. Dyna, an integrated architecture for learning, planning, and reacting. <i>ACM Sigart Bulletin</i> , 2(4):160–163, 1991.
710 711	Edan Toledo and Amanda Prorok. Codreamer: Communication-based decentralised world models. <i>arXiv preprint arXiv:2406.13600</i> , 2024.
712 713 714	Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning. In Advances in Neural Information Processing Systems, 2017. ISBN 9781510860964.
715 716 717 718	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In <i>Advances in Neural Information Processing Systems</i> , 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf.
719 720 721	Jianhao Wang, Zhizhou Ren, Terry Liu, Yang Yu, and Chongjie Zhang. QPLEX: Duplex dueling multi-agent q-learning. In <i>International Conference on Learning Representations</i> , 2021. URL https://openreview.net/forum?id=Rcmk0xxIQV.
723 724 725	Daniël Willemsen, Mario Coppola, and Guido CHE de Croon. Mambpo: Sample-efficient multi-robot reinforcement learning using learned world models. In 2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). IEEE, 2021.
726 727 728 729	Zhiwei Xu, Dapeng Li, Bin Zhang, Yuan Zhan, Yunpeng Baiia, and Guoliang Fan. Mingling foresight with imagination: Model-based cooperative multi-agent reinforcement learning. In <i>Advances in</i> <i>Neural Information Processing Systems</i> , 2022. URL https://openreview.net/forum? id=flBYpZkW6ST.
730 731 732 733	Chao Yu, Akash Velu, Eugene Vinitsky, Jiaxuan Gao, Yu Wang, Alexandre Bayen, and Yi Wu. The surprising effectiveness of PPO in cooperative multi-agent games. In <i>Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track</i> , 2022.
734 735 736	Qiaosheng Zhang, Chenjia Bai, Shuyue Hu, Zhen Wang, and Xuelong Li. Provably efficient information-directed sampling algorithms for multi-agent reinforcement learning. <i>arXiv preprint arXiv:2404.19292</i> , 2024a.
737 738 739 740	Yang Zhang, Shixin Yang, Chenjia Bai, Fei Wu, Xiu Li, Xuelong Li, and Zhen Wang. Towards efficient llm grounding for embodied multi-agent collaboration. <i>arXiv preprint arXiv:2405.14314</i> , 2024b.
741 742 743 744 745 746 747 748 749 750 751	Łukasz Kaiser, Mohammad Babaeizadeh, Piotr Miłos, Błażej Osiński, Roy H Campbell, Konrad Czechowski, Dumitru Erhan, Chelsea Finn, Piotr Kozakowski, Sergey Levine, Afroz Mohiuddin, Ryan Sepassi, George Tucker, and Henryk Michalewski. Model based reinforcement learning for atari. In <i>International Conference on Learning Representations</i> , 2020. URL https:// openreview.net/forum?id=S1xCPJHtDB.
752 753 754 755	

⁷⁵⁶ A WORLD MODELS DETAILS AND HYPERPARAMETERS

758 A.1 OBSERVATION TOKENIZER 759

Our tokenizer for local observation discretization is based on the implementation¹ of a vanilla VQ-VAE (van den Oord et al., 2017). Faced with continuous non-vision observation, we build the encoder and decoder as Multi-Layer Perceptrons (MLPs). The decoder is designed with the same hyperparameters as the ones of the encoder. The hyperparameters are listed as Table 2. During the phase of collecting experience from the external environment, each agent takes the reconstructed observations processed by the VQ-VAE as input instead to avoid the distribution shift between policy learning and policy execution.

For training this vanilla VQ-VAE, we use a straight-through estimator to enable gradient backpropa gation through the non-differentiable quantization operation in the quantization of VQ-VAE. The loss
 function for learning the autoencoder is as follows:

$$\mathcal{L}_{\text{VQ-VAE}}(E, D, \mathcal{Z}) = \mathbb{E}_{i \sim \mathcal{N}} \mathbb{E}_{o^i} \left[\| o^i - \hat{o}^i \|^2 + \| \text{sg}[E(o^i)] - z_q^i \|^2 + \beta \| \text{sg}[z_q^i] - E(o^i) \|^2 \right]$$
(7)

where $\mathcal{N} = \{1, 2, ..., n\}$ denotes the set of agents, sg[·] denotes the stop-gradient operation and β is the coefficient of the commitment loss $\|\text{sg}[z_q^i] - E(o^i)\|^2$. In practice, we found the codebook \mathcal{Z} can suffer from codebook collapse when learning from scratch. Thus, we adopt the Exponential Moving Averages (EMA) (van den Oord et al., 2017) technique to alleviate this problem.

,	-	7	1	ċ	
	1	ľ			1

770 771

110		
777	Table 2: VQV	AE hyperparameters.
778	Hyperparameter	Value
779		
780	Encoder&Decoder	
700	Layers	3
781	Hidden size	512
782		CELU(II on denselso & Cimeral 2016)
700	Activation	GELU(Hendrycks & Gimpel, 2016)
103	Codebook	
784	$\overline{\text{Codebook size } (N)}$	512
785	Tokens per observation (K)	16
786	Code dimension	128
787	Coef. of commitment loss (β)	10.0
700		

A.2 TRANSFORMER

791 The shared Transformer serving as the local dynamics model is based on the implementation of 792 minGPT (Karpathy, 2020). Given a fixed imagination horizon H, it first takes a token sequence 793 of length H(K + 1) composed of observation tokens and action tokens, and embeds it into a 794 $H(K+1) \times D$ tensor via separate embedding tables for observations and actions. Then, the aggregated feature tensor, returned by the agent-wise aggregation module, is inserted after the action 796 embedding tensor at every timestep, forming a final embedding tensor of shape $H(K+2) \times D$. This 797 tensor is forwarded through fixed Transformer blocks. Here, we adopt GPT2-like blocks (Radford et al., 2019) as the basic blocks. The hyperparameters are listed as Table 3. To enable training across 798 all environments on a single NVIDIA RTX 3090 GPU, we adapt imagination horizon H based on the 799 number of agents. 800

801 802

789

790

A.3 PERCEIVER

The Perceiver (Jaegle et al., 2021) is based on the open-source implementation². By aligning the length of the latent querying array with the number of agents n, we obtain the intrinsic global representation feature corresponding to each individual agent. We further dive into the process of agent-wise representation aggregation: (i) the embedding tensor of shape $(K + 1) \times D$ at each timestep, mentioned in Appendix A.2, is concatenated with others from all agents, thereby getting a

¹Code can be found in https://github.com/lucidrains/vector-quantize-pytorch ²Code can be found in https://github.com/lucidrains/perceiver-pytorch

811	Table 3: Transformer hyp	erparameters.
812	Hyperparameter	Value
813	Imagination horizon (H)	{15 8 5}
814	Embedding dimension	256
815	Lavers	10
816	Attention heads	4
817	Weight decay	0.01
818	Embedding dropout	0.1
819	Attention dropout	0.1
820	Residual dropout	0.1
821		
822		
823	Table 4: Perceiver hyper	rparameters.
824	Hyperparameter	Value
825	Length of latent querying	n (number of agents
826	Cross attention heads	8
827	Inner Transformer layers	2
828	Transformer attention heads	8
829	Dimension per attention head	64
830	Embedding dropout	0.1
831	Attention dropout	0.1
832	Residual dropout	0.1

 $n(K+1) \times D$ sequence for the joint observation-action pair at the current timestep; (ii) through the cross-attention mechanism with the latent querying array, the original sequence is compressed from length n(K+1) to n; (iii) the compressed sequence is then forwarded through a standard transformer with bidirectional attention inside the Perceiver. The hyperparameters are listed as Table 4.

B BEHAVIOUR LEARNING DETAILS

In MARIE, we use MAPPO-like (Yu et al., 2022) actor and critic, where the actor and critic should have been 3-layer MLPs. However, unlike other CTDE model-free approaches, whose critic takes additional global oracle states from the environment in the training phase, our world model hardly provides related predictions in the imagined trajectories. To alleviate this issue, we augment the critic with an attention mechanism and provide it all reconstructed observations \hat{o}_t of all agents. Therefore, the actor ψ remains a 3-layer MLP with ReLU activation, while the critic ξ is enhanced with an extra layer of self-attention, built on top of the original 3-layer MLP, i.e., we overwrite the critic $V_{\xi}^{i}(\hat{\boldsymbol{o}}_{t}) \simeq \mathbb{E}_{\pi_{\psi}^{i}}(\sum_{l \geq t} \gamma^{l-t} \hat{r}_{l}^{i})$ for agent *i*. Similar to off-the-shelf CTDE model-free approaches, we adopt parameter sharing across agents.

Critic loss function We utilize λ -return in Dreamer (Hafner et al., 2020), which employs an 852 exponentially-weighted average of different k-steps TD targets to balance bias and variance as the 853 regression target for the critic. Given an imagined trajectory $\{\hat{o}_{\tau}^{i}, a_{\tau}^{i}, \hat{r}_{\tau}^{i}, \hat{\gamma}_{\tau}^{i}\}_{t=1}^{H}$ for agent i, λ -return 854 is calculated recursively as,

$$V_{\lambda}^{i}(\hat{\boldsymbol{o}}_{t}) = \begin{cases} \hat{r}_{t}^{i} + \hat{\gamma}_{t}^{i} \left[(1-\lambda) V_{\xi}^{i}(\hat{\boldsymbol{o}}_{t}) + \lambda V_{\lambda}^{i}(\hat{\boldsymbol{o}}_{t+1}) \right] & \text{if } t < H \\ V_{\xi}^{i}(\hat{\boldsymbol{o}}_{t}) & \text{if } t = H \end{cases}$$

$$\tag{8}$$

The objective of the critic ξ is to minimize the mean squared difference \mathcal{L}^i_{ξ} with λ -returns over imagined trajectories for each agent *i*, as

$$\mathcal{L}_{\xi}^{i} = \mathbb{E}_{\pi_{\psi}^{i}} \left[\sum_{t=1}^{H-1} \left(V_{\xi}^{i}(\hat{\boldsymbol{o}}_{t}) - \operatorname{sg}(V_{\lambda}^{i}(\hat{\boldsymbol{o}}_{t})) \right)^{2} \right]$$
(9)

where sg(·) denotes the stop-gradient operation. We optimize the critic loss with respect to the critic parameters ξ using the Adam optimizer.

Hyperparameter	Value
Imagination Horizon (H)	{15, 8, 5
Predicted discount label γ	Ò.99
λ	0.95
η	0.001
Clipping parameter ϵ	0.2

873

864

874 Actor loss function The objective for the action model $\pi_{\psi}(\cdot | \hat{o}_t^i)$ is to output actions that maximize 875 the prediction of long-term future rewards made by the critic. To incorporate intermediate rewards 876 more directly, we train the actor to maximize the same λ -return that was computed for training the 877 critic. In terms of the non-stationarity issue in multi-agent scenarios, we adopt PPO updates, which 878 introduce important sampling for actor learning. The actor loss function for agent *i* is:

$$\mathcal{L}^{i}_{\psi} = -\mathbb{E}_{p_{\phi}, \pi^{i}_{\psi_{\text{old}}}} \left[\sum_{t=0}^{H-1} \min\left(r^{i}_{t}(\psi) A^{i}_{t}, \operatorname{clip}(r^{i}_{t}(\psi), 1-\epsilon, 1+\epsilon) A^{i}_{t} \right) + \eta \mathcal{H}(\pi^{i}_{\psi}(\cdot|\hat{o}^{i}_{t})) \right]$$
(10)

where $r_t^i(\psi) = \pi_{\psi}^i/\pi_{\psi_{\text{old}}}^i$ is the policy ratio and $A_t^i = \operatorname{sg}(V_{\lambda}^i(\hat{o}_t) - V_{\xi}^i(\hat{o}_t))$ is the advantage. We optimize the actor loss with respect to the actor parameters ψ using the Adam optimizer. In the discount prediction of MARIE, we set its learning target γ to be 0.99. Overall hyperparameters are shown in Table 5.

904

906

883

884

C EXTENDED ANALYSIS ON ATTENTION PATTERNS

To provide qualitative analysis of our world model, we select typical scenarios $-3s_vs_5z$ where our method achieves the most significant improvement compared to other baselines for visualizing attention maps inside the Transformer. For the sake of simple and clear visualization, we set the imagination horizon *H* as 5. In terms of cross-attention maps in the aggregation module, we select a scenario 2s3z including 5 agents for visualization. Visualization results are depicted as Fig. 8 and Fig. 9.

895 The prediction of local dynamics entails two distinct attention patterns. The left one in Fig. 8 can 896 be interpreted as a Markovian pattern, in which the observation prediction lays its focus on the 897 previous transition. In contrast, the right one is regularly striated, with the model attending to specific tokens in multiple prior observations. In terms of the agent-wise aggregation, we also identify two 899 distinct patterns: *individuality* and *commonality*. The top one in Fig. 9 illustrates that each agent flexibly attends to different tokens according to their specific needs. In contrast, the bottom one 900 exhibits consistent attention allocation across all agents, with attention highlighted in nearly identical 901 positions. The diverse patterns in the Transformer and Perceiver may be the key to accurate and 902 consistent imagination. 903

905 D BASELINE IMPLEMENTATION DETAILS

MAMBA (Egorov & Shpilman, 2022) is evaluated based on the open-source implementation:
 https://github.com/jbr-ai-labs/mamba with the hyperparameters in Table 6.

909 MAPPO (Yu et al., 2022) is evaluated based on the open-source implementation: https://github.com/marlbenchmark/on-policy with the common hyperparameters in Table 7.
 911

912 QMIX (Rashid et al., 2018) is evaluated based on the open-source implementation: https://github.com/oxwhirl/pymarl with the hyperparameters in Table 8.

914 QPLEX (Wang et al., 2021) is evaluated based on the open-source implementation: https:
 915 //github.com/wjh720/QPLEX with the hyperparameters in Table 9. Since its implementation is mostly based on the open-source implementation: PyMARL (Samvelyan et al., 2019), its most hyperparameters setting remains the same as the one in QMIX in addition to its own special hyperparameters.

921			
922	Table 6: Hyperparameters f	for MAMBA in SM.	AC environments.
923	Hyperparameter		Value
924	Batch size		256
925	λ for λ -return comp	utation	0.95
926	Entropy coefficient		0.001
927	Entropy annealing		0.99998
928	Number of policy up	odates	4
929	Epochs per policy up	odate	5
930	Clipping parameter of	ŝ	0.2
931	Actor Learning rate		0.0005
932	Critic Learning rate		0.0005
933	Discount factor γ		0.99
934	Model Learning rate		0.0002
935	Number of model tra	aining epochs	60
006	Number of imagined	l rollouts	800
930	Sequence length		20
937	Imagination horizon	H	15
938	Buffer size	_	2.5×10^{5}
939	Number of categoric	als	32
940	Number of classes		32
941	KL balancing entrop	y weight	0.2
942	KL balancing cross of	entropy weight	0.8
943	Gradient clipping	. 1	100
944	Collected trajectorie	s between updates	1
945	Hidden size		256
946			
947			
948			
949			
950			
951			
952			
953	Table 7: Common hyperparame	eters for MAPPO in	SMAC environments.
954	Hyperparameter	Value	
955 -		1 00	
956	Batch size	num envs \times buffer	length \times num agents
957	Mini batch size	batch size / mini-b	atch
958	Recurrent data chunk length	10	
050	$GAE \lambda$	0.95	
909	Discount factor γ	0.99 hashan lana	
960	Value loss	nuber loss	
961	Auber della	10.0 A dam	
962	Optimizer learning rate	Auani 0.0005	
963	Optimizer ensilon	$1 \sim 10^{-5}$	
964	Weight decay	$1 \land 10$	
965	Gradient clipping	10	
966	Network initialization	orthogonal	
967	Use reward normalization	True	
968	Use feature normalization	True	
969 -	est reature normanization		
970			

972			
973			
974			
975			
976			
977			
978			
979	Table 8: Huperpersonators for OMIX is	n SMAC anvir	onmonto
980	Table 8. Hyperparameters for QMIX II		onnents.
981	Hyperparameter	Value	
982	Batch size	32	
983	Buffer size	5000	
984	Epsilon in epsilon-greedy	$1.0 \rightarrow 0.05$	
985	Epsilon anneal time	50000	
986	Train interval	1 episode	
987	Discount factor γ	0.99	
988	Optimizer DMSDrop or	RMSProp	
989	\mathbf{R}	0.99 10^{-5}	
990	Gradient clipping	10	
991		10	
992			
993			
994			
995			
996			
997			
998			
999			
1000			
1001			
1002			
1003			
1004			
1005			
1006	Table 9: Hyperparameters for OPLEX i	n SMAC envi	ronments.
1007	Hypernarameter	Vəluc	
1008		Value	
1009	Batch size	32	
1010	Buffer size	5000	0.05
1011	Epsilon in epsilon-greedy	1.0 - 50000	→ 0.05)
1012	Epshon annear time	30000 1 ania) Jada
1013	Discount factor γ		oue
1014	Ontimizer	RMS	Pron
1015	RMSProp α	0.99	Top
1016	RMSProp ϵ	10^{-5}	
1017	Gradient clipping	10	
1018	Number of layers in HyperNetwork	. 1	
1019	Number of heads in the attention me	odule 4	
1020			
1021			
1022			
1023			
1024			
1025			



Figure 8: Attention patterns in the Transformer. We observe two distinct types of attention weights during the prediction of local dynamics. In the first one (*left*), the next observation prediction is primarily dependent on the last transition, which means the world model has learned the Markov property corresponding to Dec-POMDPs. The second type (*right*) exhibits a regularly striated pattern, where the next observation prediction hinges overwhelmingly on the same dimension of multiple previous timesteps. The above attention weights are produced by a sixth-layer and ninth-layer attention head during imaginations on the $3s_vs_s5z$ scenario.



Figure 9: Cross attention patterns in the Perceiver. We observe the *individuality* and *commonality* in the agent-wise aggregation. The top part of the figure represents the *individuality*, where agents adjust their attentions over the whole joint token sequence at timestep t flexibly according to their own needs. In contrast, the bottom exhibits the *commonality*, where every agent's attention over the joint token sequence is emphasized in the similar positions of the sequence. The cross attention weights mentioned above are produced by the first and sixth head of the cross attention within the Perceiver, during the agent-wise aggregation on the 2s3z scenario.

MBVD (Xu et al., 2022) is evaluated based on the implementation in its supplementary material from https://openreview.net/forum?id=flBYpZkW6ST with the hyperparameters in Table 10. Akin to QPLEX, its implementation is based on the open-source implementation: PyMARL, its most hyperparameters setting remains the same as the one in QMIX in addition to its own special hyperparameters.

1069

1063

1041

1042

1043

1044

1045

1046

1047

1048 1049 1050

1051

1052 1053 1054

1055

1070 E ADDITIONAL EXPERIMENTS

1072

E.1 EVALUATIONS ON MAMUJOCO

The Multi-Agent MuJoCo (MAMuJoCo) (Peng et al., 2021) environment is a multi-agent extension of
 MuJoCo. While the MuJoCo tasks challenge a robot to learn an optimal way of motion, MAMuJoCo
 models each part of a robot as an independent agent — for example, a leg for a spider or an
 arm for a swimmer — and requires the agents to collectively perform efficient motion. With the
 increasing variety of the body parts, MAMujoco can be also considered as a testbed for evaluating the
 coordination among heterogeneous agents, which poses a big challenge for learning the multi-agent
 dynamics inside it, especially in a *decentralized* manner.

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111 1112

082	Hyperparameter	Value
083	Batch size	32
)84	Buffer size	5000
)85	Epsilon in epsilon-greedy	$1.0 \rightarrow 0.05$
)86	Epsilon anneal time	50000
)87	Train interval	1 episode
88	Discount factor γ	0.99
)89	Optimizer	RMSProp
)90	RMSProp α	0.99
91	RMSProp ϵ	10^{-5}
92	Gradient clipping	10
93	Number of layers in HyperNetwork	1
94	Number of heads in the attention module	4
95	Horizon of the imagined rollout	3
96	KL balancing α	0.3
90	Dimension of the latent state \hat{s}	num agents x 16
191		
190		MARIE
199		
100	HalfCheetah-v2-2x3 HalfCheetah-v2-3x2	Wal



Figure 10: Curves of the performance for MARIE, MAMBA, HAPPO and MAPPO in 3 chosen MAMujoco scenarios. Y axis: return; X axis: number of steps taken in the real environment.

1113 While MAMBA (Egorov & Shpilman, 2022) originally does not take the continuous action space 1114 case into consideration, which is a obvious limitation in it, we would like to evaluate MARIE in such case, e.g., MAMujoco, to better demonstrate that our method can be also effective in other 1115 multi-agent domains. In MAMujoco, we discretize the scalar in each dimension of continuous actions 1116 into one of 256 fixed-width bins independently to obtain discrete action tokens for local dynamics 1117 learning. As the behaviour learning in MARIE adopts a MAPPO-like and on-policy manner, we 1118 choose two strong on-policy PPO-based baselines - MAPPO (Yu et al., 2022) and HAPPO (Kuba 1119 et al., 2021). Additionally, we also include MAMBA as a model-based baseline for comparison. Since 1120 MAMBA(Egorov & Shpilman, 2022) was originally originally designed for domains with discrete 1121 action space, significant effort was required to adapt and evaluate it on MAMujoco, which features 1122 continuous action space. The experiments are conducted in HalfCheetah-v2-2x3, HalfCheetah-v2-3x2 1123 and Walker2d-v2-2x3. The learning curves of the return averaged over 4 seeds are presented as 1124 Figure 10. Notably, MAMBA fails to enhance policy learning in the Walker2d-v2-2x3 scenario and remains exceptionally time-consuming. Consequently, we report its results only for 1 million 1125 environment steps in this scenario. 1126

As illustrated in Figure 10, our MARIE consistently shows superior sample efficiency and achieves the best performance in 2 of 3 scenarios with limited 2M environment steps. For the performance difference between MAMBA and MARIE in *HalfCheetah-v2-3x2*, we hypothesize that MAMBA's policy learning benefits significantly from using the internal recurrent features of the world model as inputs in this scenario, while the policy in our method only takes the reconstructed observation as input in order to support fast deployment in the environment without the participation of the world model. We attribute the performance gap between MARIE and other two model-free baselines in *HalfCheetah-v2-3x2* to the access to global oracle state in the chosen baselines. The policy in our



Figure 11: Curves of the performance for MARIE, MAMBA, QMIX, HAPPO and MAPPO in 3 chosen SMACv2 scenarios. Y axis: Win Rate; X axis: number of steps taken in the real environment. While MARIE shows competitive performance to MAMBA on zerg_5_vs_5, MARIE is superior to all other baselines in terms of sample efficiency and final performance in the rest 2 scenarios.

algorithm is purely learned from the inner imaginations of the world model where there is only
 reconstructed local observation. Considering MAMujoco is a multi-heterogeneous-agent benchmark
 which necessitates a more precise credit assignment during training, it would be much more helpful
 for policy learning to have access to the true global oracle state than in other benchmarks. But overall,
 our MARIE presents a faster convergence rate, implying that our Transformer-based world model
 can generate accurate imaginations and bring remarkable sample efficiency.

1155 1156 1157

1149

E.2 EVALUATIONS ON SMACV2

1158 Given known serious flaws in SMACv1 (e.g., the tricky open-loop policy issue), we extend our 1159 evaluation of MARIE to SMACv2 (Ellis et al., 2023), which introduces more stochasticity and partial 1160 observability. In this comparison, we benchmark MARIE against four baselines: MAPPO, HAPPO, 1161 QMIX and MAMBA. For each random seed, we adopt the same evaluation protocol as the main experiment on SMACv1. Importantly, the hyperparameters of MARIE remain unchanged, as detailed 1162 in §H. Here, we directly use the results of MAPPO and QMIX provided in the official SMACv2 1163 repository³. Illustrated in Figure 11, while MARIE shows competitive performance to MAMBA 1164 on zerg_5_vs_5, MARIE is superior to all other baselines in terms of sample efficiency and final 1165 performance in the rest 2 scenarios. 1166

- 1167
- 1168

E.3 COMPARISON WITH QMIX AND QPLEX WITH DIFFERENT EPSILON ANNEALING TIME

1169 Considering the potentially inappropriate influence of a large ϵ annealing time used in the epsilon-1170 greedy algorithm when evaluated in the low data regime evaluation, we run QMIX and QPLEX 1171 with a smaller ϵ annealing time, and compare the performance of them with ours. The result is 1172 reported in Table 11. The reported result shows that the original hyperparameters used in the main 1173 experiment, which are also directly referred to Egorov & Shpilman (2022) and Liu et al. (2024), 1174 are reasonable since the performance of QMIX and QPLEX under the original hyperparameters is superior to the ones with a smaller ϵ annealing time at most scenarios. Besides, our MARIE still 1175 consistently outperforms QMIX and QPLEX with a smaller ϵ annealing time. 1176

1177

1178 E.4 COMPARISON WITH EXISTING TRANSFORMER-BASED WORLD MODELS

Existing Transformer-based world models are primarily designed for single-agent scenarios, but they can be naturally adapted to multi-agent settings, modeling either independently local dynamics or joint dynamics. Fortunately, we have included IRIS as a Transformer-based world model baseline in our ablation experiments. Specifically, the *Centralized Manner* and *MARIE w/o aggregation* variants from our ablation experiments correspond to IRIS baseline variants under different deployment strategies. But different from their original implementation, these IRIS baseline variants also uses the same actor-critic method as MARIE during learning in imaginations phase (i.e., using PPO instead

³**Results of QMIX and MAPPO are available at** https://github.com/oxwhirl/smacv2/tree/ main/smacv2/examples/results.

1189	Table 11: Mean evaluation win rate and standard deviation for QMIX and QPLEX with different
1190	epsilon anneal time t_{ϵ} over 4 random seeds. We bold the values of the maximum.

191							
1192	Maps	Steps	MARIE	QMIX ($t_{\epsilon} = 50000$)	QMIX ($t_{\epsilon} = 10000$)	QPLEX ($t_{\epsilon} = 50000$)	QPLEX ($t_{\epsilon} = 10000$)
1193	1c3s5z		85.0 (9.4)	43.6(29.2)	33.3(15.0)	68.3(7.4)	44.8(11.0)
110/	2m_vs_1z		95.5 (7.9)	70.3(14.8)	36.1(28.2)	84.8(10.8)	93.2(4.7)
1194	2s_vs_1sc		96.9 (7.1)	0.0(0.0)	3.9(6.7)	15.7(19.5)	43.2(32.4)
1195	2s3z	10012	80.5(9.3)	37.7(15.5)	29.1(20.3)	50.2(8.4)	28.3(11.5)
1196	3m	100K	99.5 (0.4)	54.4(22.7)	63.8(14.6)	88.7(6.9)	85.0(11.3)
1197	3s_vs_3z		98.9 (1.5)	0.0(0.0)	0.0(0.0)	0.0(0.0)	0.0(0.0)
1198	3s_vs_4z		73.0 (6.2)	0.0(0.0)	0.0(0.0)	0.0(0.0)	0.0(0.0)
1100	8m		88.0 (3.9)	69.5(12.8)	68.6(13.6)	83.4(6.4)	79.7(9.8)
1000	MMM		87.6 (3.0)	31.1(17.3)	18.9(4.3)	69.3(35.1)	20.2(7.7)
1200	so_many_baneling		94.8 (5.9)	20.0(8.9)	30.7(18.5)	32.2(6.1)	37.7(9.2)
1201	2c_vs_64zg	2001	25.9 (14.3)	0.5(0.5)	0.0(0.0)	0.1(0.1)	0.0(0.0)
1202	3s_vs_5z	200 K	78.4 (11.2)	0.0(0.0)	0.0(0.0)	0.0(0.0)	0.0(0.0)
1203 1204	corridor	400K	71.0 (13.8)	0.0(0.0)	0.0(0.0)	0.0(0.0)	0.0(0.0)



Figure 12: Comparison with two direct extension of IRIS in the *so_many_baneling* scenario.

of REINFORCE for behaviour learning). With a shared behaviour learning phase, we can analyze the reason why existing single-agent Transformer-based world model cannot be directly adapted to MARL. As shown in Figure 12, without incorporating CTDE principle, the learning of single-agent world model would be disrupted by the scalability and non-stationarity issues.

E.5 COMPARISON AGAINST MAMBA WITH DIFFERENT IMAGINATION HORIZON

We report the performance of MARIE and MAMBA on $3s_vs_sz_v$ with different imagination horizon in Table 12. And the result shows that the performance gap between the two is not related to the choice of imagination horizon. Interestingly, a larger imagination horizon may help policy learning in imagination. We hypothesize that longer imagined trajectories help alleviating shortsighted behaviours in policy learning.

Table 12: The performance of MARIE and MAMBA on $3s_{vs}5z$ with different imagination horizon.

1238 1239	Maps	Method	Horizon $H = 8$	Horizon $H = 15$	Horizon $H = 25$
1240 1241	3s_vs_5z	MARIE MAMBA	$0.40 \pm 0.34 \\ 0.00 \pm 0.00$	$0.75 \pm 0.09 \\ 0.13 \pm 0.14$	0.78 ± 0.11 0.16 ± 0.13

F ADDITIONAL DISCUSSION BETWEEN CODREAMER AND MARIE

Additionally, a recent method CoDreamer (Toledo & Prorok, 2024) extends DreamerV3 (Hafner et al., 2023) to the multi-agent setting, using GAT V2 (Brody et al., 2021) for communication among agents' world models and policies. Though the aggregation modules in CoDreamer and ours are both built upon the Transformer architecture, our focus lies in computational efficiency of aggregation while it focuses on the underlying topological graph structure among agents. However, a fundamental difference is the backbone used for modeling the local dynamics. While we cast the local dynamics learning as the sequence modeling over discrete tokens, which can be achieved by using auto-regressive Transformers with causal attention mechanism, CoDreamer directly adopts the RSSM framework in DreamerV3.

STANDARDIZED PERFORMANCE EVALUATION PROTOCOL G

Agarwal et al. (2021) discuss the limitations of mean and median scores, and show that substantial discrepancies arise between standard point estimates and interval estimates in RL benchmarks. To deliver a rigorous statistical evaluation, we summarize in Figure 13 the win rate with stratified bootstrap confidence intervals for mean, median, and inter-quartile mean (IQM). For finer comparisons, we also provide probabilities of improvement in Figure 14.



Figure 13: Mean, median, and inter-quartile mean win rate, computed with stratified bootstrap confidence intervals. 4 runs for all algorithms.



Figure 14: Probabilities of improvement (Agarwal et al., 2021).

Hyperparameters Hyperparameter Batch size for tokenizer training Batch size for world model training Optimizer for tokenizer	Value 256
Hyperparameter 299 Batch size for tokenizer training 300 Batch size for world model training 301 Optimizer for tokenizer	256
299Batch size for tokenizer training300Batch size for world model training301Optimizer for tokenizer	256
Batch size for world model trainingOptimizer for tokenizer	20
301 Optimizer for tokenizer	30
	AdamW
302 Optimizer for world model	AdamW
303 Optimizer for actor & critic	Adam
304 Tokenizer learning rate	0.0003
World model learning rate	0.0001
Actor learning rate	0.0005
307 Critic learning rate	0.0005
Gradient clipping for actor & critic	100
Gradient clipping for tokenizer	10
Gradient clipping for world model	10
Weight decay for world model	0.01
λ for λ -return computation	0.95
Discount factor γ	0.99
Entropy coefficient	0.001
Buffer size (transitions)	$2.5 imes 10^5$
Number of tokenizer training epochs	200
16 Number of world model training epochs	200
17 Collected <u>transitions</u> between updates	$\{100, 200\}$
Epochs per policy update (PPO epochs)	5
PPO Clipping parameter ϵ	0.2
Number of imagined rollouts	600 or 400
Imagination horizon H	$\{15, 8, 5\}$
Number of policy updates	$\{4, 10, 30\}$
Number of stacking observations	5
Observe agent id	False
Observe last action of itself	False
326	
327	

Table 14: Computational time consumption of MARIE in SMAC.

Environment Steps	100000	200000	400000
Training Time	1 day	2-3 days	4 days
Usage of GPU Mem	22GB	22GB	22GB

H PARAMETERS SETTING AND COMPUTATIONAL CONSUMPTION OF MARIE

All our experiments are run on a machine with a single NVIDIA RTX 3090 GPU, a 36-core CPU, and 128GB RAM. We provide the hyperparameters of MARIE for experiments in SMAC, shown as Table 13. To enable the running of experiments in all SMAC scenarios with a single NVIDIA RTX 3090 GPU, we set the imagination horizon H as 8 for other scenarios involving the number of agents n > 5, 15 for $n \le 5$. In so_many_baneling and 2s3z, we set the imagination horizon H as 5. Correspondingly, the number of policy updates in imaginations varies with imagination horizon H. As for the scenario $2c_{vs}_{ds}$, considering the significantly large action space in it, we enable the observation of agent id and last action for each agent and disable stacking the last 5 observations as input to the policy.

Based on the above reported setting, we present a rough computational consumption in Table 14.

1350 I OVERVIEW OF MARIE ALGORITHM 1351

1352 Pseudo-code is summarized as Algorithm 1.

٩lg	orithm 1 MARIE
//	main loop of training
f	or epochs do
	collect_experience(num_transitions)
	for learning_world_model_steps_per_epoch do
	train_world_model()
	end for
	for learning_behaviour_steps_per_epoch do
	train_agents()
	end for
e	nd for
f	unction collect experience (n) :
0	$\phi \leftarrow env reset()$
f	arr t = 0 $n-1$ do
1	// processed by VO-VAE
	$\hat{\boldsymbol{o}} \leftarrow D(E(\boldsymbol{o}))$
	Sample $a^i \sim \pi^i (a^i \hat{a}^i) \forall i$
	$\mathbf{o}_{t+1} r_t done \leftarrow \text{env sten}(\mathbf{a}_t)$
	if done $-True$ then
	$\mathbf{O}_{i+1} \leftarrow \text{env} \text{ reset}(\mathbf{O}_{i+1})$
	$v_{t+1} \leftarrow 0$
	lse
	$\gamma_t \leftarrow 0.99$
	end if
е	nd for
Ĩ	$\mathcal{D} \leftarrow \mathcal{D} \cup \{ \boldsymbol{\alpha}_{t} \boldsymbol{\alpha}_{t} \boldsymbol{\gamma}_{t} \rangle \}_{t=0}^{n-1}$
-	$(2) = (0, 1, \infty, 1, 1, 1)_{t=0}$
f	unction train_world_model():
S	ample $\{\boldsymbol{o}_t, \boldsymbol{a}_t, r_t, \gamma_t\}_{t=\tau}^{t=\tau+H-1}$
ι	Update (E, D, \mathcal{Z}) via \mathcal{L}_{VQ-VAE} over observations $\{o_t\}_{t=\tau}^{t=\tau+H-1}$
f	for agent $i = 1, \ldots, n$ do
	Update ϕ, θ via $\mathcal{L}_{\text{Dyn}}(\phi, \theta)$ over local trajectories $\{o_i^i, a_t^i, r_t, \gamma_t\}_{t=\tau}^{t=\tau+H-1}$
e	nd for
f	unction train_agents():
S	ample an initial observation $o_0 \sim D$
{	$x_{0}^{i}]_{i=1}^{K} \leftarrow E(o_{0}^{i}), \hat{o}_{0}^{i} \leftarrow D(E(o_{0}^{i})), \forall i$
f	or $t = 0,, H - 1$ do
	Sample $a_t^i \sim \pi_{\psi}^i(a_t^i \hat{o}_t^i), \forall i$
	Aggregate $(x_{t_1}^1, \ldots, x_{t_K}^n, a_{t_1}^1, \ldots, x_{t_1}^n, \ldots, x_{t_K}^n, a_t^n)$ into (e_1^1, \ldots, e_t^n) via the Perceiver θ
	Sample \hat{r}^i
	$\hat{a}_{i} = D(\hat{a}_{i}) \forall i$
-	$O_{t+1} \leftarrow D(x_{t+1,\cdot}), \forall t$
e	nu ior
I	or agent $i = 1, \dots, n$ do
	Update actor π_{ψ}^{i} and critic V_{ξ}^{i} via $\mathcal{L}_{\text{Dyn}}(\phi, \theta)$ over imagined trajectories $\{\hat{o}_{t}^{i}, a_{t}^{i}, \hat{r}_{t}^{i}, \hat{\gamma}_{t}^{i}\}_{t=0}^{t=H}$
~	nd for

- 1402
- 1403