Efficient Information Sharing for Training Decentralized Multi-Agent World Models

Xiaoling Zeng, Qi Zhang

Keywords: cooperative multi-agent reinforcement learning, decentralized world models, multi-agent communication.

Summary

World models, which were originally developed for single-agent reinforcement learning, have recently been extended to multi-agent settings. Due to unique challenges in multi-agent reinforcement learning, agents' independent training of their world models often leads to underperforming policies, and therefore existing work has largely been limited to the centralized training framework that requires excessive communication. As communication is key, we ask the question of how the agents should communicate efficiently to train and learn policies from their decentralized world models. We address this question progressively. We first allow the agents to communicate with unlimited bandwidth to identify which algorithmic components would benefit the most from what types of communication. Then, we restrict the inter-agent communication with a predetermined bandwidth limit to challenge the agents to communicate efficiently. Our algorithmic innovations develop a scheme that prioritizes important information to share while respecting the bandwidth limit. The resulting method yields superior sample efficiency, sometimes even over centralized training baselines, in a range of cooperative multiagent reinforcement learning benchmarks.

Contribution(s)

- 1. This paper proposes a model-based MARL method that explicitly considers communication in both the world model and the actor-critic training stages, and analyzes the impact of communication bandwidth on decentralized training.
 - **Context:** Previous work either builds on model-free MARL, discussing only experience sharing under different bandwidths with limited model and information diversity, or extracts shared agent information as features for centralized training, but omits these features during decentralized execution (Gerstgrasser et al., 2023; Venugopal et al., 2023).
- Our work systematically studies information sharing under bandwidth limitations, aiming to optimize communication efficiency while guaranteeing performance within a decentralized framework.
 - **Context:** Existing methods fail to effectively deal with communication bandwidth constraints and often rely on Euclidean distance constraints to filter communication neighbors (Toledo & Prorok, 2024).

Efficient Information Sharing for Training Decentralized Multi-Agent World Models

Xiaoling Zeng¹, Qi Zhang^{1,†}

xz13@email.sc.edu, qzhang9@wpi.edu

¹Department of Computer Science and Engineering, University of South Carolina

²Department of Computer Science, Worcester Polytechnic Institute

[†] Work done when employed at the University of South Carolina

Abstract

World models, which were originally developed for single-agent reinforcement learning, have recently been extended to multi-agent settings. Due to unique challenges in multi-agent reinforcement learning, agents' independent training of their world models often leads to underperforming policies, and therefore existing work has largely been limited to the centralized training framework that requires excessive communication. As communication is key, we ask the question of how the agents should communicate efficiently to train and learn policies from their decentralized world models. We address this question progressively. We first allow the agents to communicate with unlimited bandwidth to identify which algorithmic components would benefit the most from what types of communication. Then, we restrict the inter-agent communication with a predetermined bandwidth limit to challenge the agents to communicate efficiently. Our algorithmic innovations develop a scheme that prioritizes important information to share while respecting the bandwidth limit. The resulting method yields superior sample efficiency, sometimes even over centralized training baselines, in a range of cooperative multi-agent reinforcement learning benchmarks.

1 Introduction

In model-based reinforcement learning (RL), the agent learns a world model that encodes its raw observations to latent states in a way that effectively recovers/predicts the observations, rewards, and future latent state dynamics. This model-based framework has contributed algorithms that have been shown to greatly improve sample efficiency for single-agent RL (Hafner et al., 2019b;a; Schrittwieser et al., 2020; Ye et al., 2021). In this paper, we are interested in extending the success of world models to the setting of cooperative multi-agent reinforcement learning (MARL) where a team of agents collectively interact in an environment to achieve a shared goal, which finds a wide range of applications such as video games (Vinyals et al., 2019), traffic and vehicle control (Chu et al., 2019; Dinneweth et al., 2022), and multi-robot systems (Corke et al., 2005). Such scenarios introduce additional challenges on top of single-agent RL, such as partial observability when the agents only partially observe the environment (Oliehoek et al., 2016) and non-stationarity as all agents concurrently update their policies during training, causing the environment dynamics to continuously change from the perspective of any individual agent (Hernandez-Leal et al., 2017). Due to these challenges, existing success in learning multi-agent world models has been primarily achieved through centralized training, where a single world model is trained and shared by all agents (Egorov & Shpilman, 2022; Venugopal et al., 2023).

Although effective, centralized training requires substantial inter-agent communication, limiting its applicability and scalability. On the other hand, as confirmed in previous work (Toledo & Prorok, 2024) and this paper, independent learning of multi-agent world models without explicit communication results in ineffective multi-agent policies after planning with the world models. This raises a critical question: *How should the agents communicate efficiently to train and learn policies from their decentralized world models?* Addressing this question is particularly challenging, since a world model, which is typically a latent-space model that captures transition dynamics, observations, and rewards, consists of inter-dependent components that involve various types of information for communication, which is further complicated by the bandwidth limitations. Adopting DreamerV2 (Hafner et al., 2020) as the architecture backbone of our decentralized world models, this work addresses this question with a two-stage study.

In the first stage, we allow agents to have unlimited communication bandwidth so that we can focus on identifying which algorithmic components would benefit the most from what types of communication. Specifically, we separately allow the information to be shared in an unlimited manner between the agents for the decentralized training of their local world models and actor-critic networks, respectively. Evaluated on the cooperative MARL benchmark SMAC, both types of information sharing yield multi-agent policies that outperform 1) the no communication baseline by a large margin and 2) the centralized training baseline in some SMAC scenarios. This might be surprising, as centralized training is widely considered a performance upper bound. We attribute this to selective communication and modular learning, which help reduce overfitting, avoid interference from mixed experiences, and support task-specific specialization. Encouraged by the results from the first stage, the second stage restricts the inter-agent communication with a predetermined bandwidth limit, which further challenges the agents to efficiently communicate with selective information. By experimenting with various bandwidths ranging from small to the largest (i.e., unlimited), our results show that there exists a relatively small bandwidth that works well for both types of information sharing, the performance of which is comparable or even better than that with unlimited bandwidth. This indicates that selective sharing under limited bandwidth can filter out noisy or redundant information, leading to more stable and effective learning.

2 Related work

Single- and multi-agent world models. One of the earliest model-based RL algorithms is Dyna (Sutton, 1991), in which the agent alternates between learning a world model of state dynamics with reward signals and planning with it to take an action. Dyna's framework has been adopted in many recent model-based RL algorithms including the Dreamer family (Hafner et al., 2019a; 2020; 2023), which is known for their effectiveness in addressing partial observability and simplicity of training a policy from the learned world model. This paper and most recent works on multi-agent world models use Dreamer as the architecture backbone. Egorov & Shpilman (2022) develop MAMBA, a centralized training and centralized execution framework where all agents share their local observations in both the global world model and the local policies. Adapting MAMBA, Venugopal et al. (2023) introduce MABL that employs a bi-level hierarchy to enhance the agents' understanding of global information during centralized world model training, while enabling fully distributed local policies for execution. Xu et al. (2022) consider model-based cooperative MARL in the centralized training framework of value decomposition. In contrast to these works, in this work, each agent maintains a local world model and learns it in a decentralized manner with inter-agent communication.

Decentralized MARL with communication. Due to challenges such as partial observability and non-stationarity, effective training of cooperative MARL agents requires either centralization like a centralized critic (Lowe et al., 2017; Chu et al., 2019; Rashid et al., 2020) and a global world model (Egorov & Shpilman, 2022; Venugopal et al., 2023) or inter-agent communication of learnable parameters (Chen et al., 2022), experiences of local trajectories (Christianos et al., 2020; Gerstgrasser et al., 2023), intents (Kim et al., 2020), etc. Closest to our work is CoDreamer (Toledo & Prorok, 2024), where the agents communicate over a graph to train a centralized world model that encodes

and averages the agents' observations and actions with a graph neural network, and therefore it falls into the centralized training framework. In contrast, our work trains decentralized world models where each agent maintains its local world model and policy. Moreover, CoDreamer operates on a predefined graph for communication, while our work focuses on achieving efficient communication where agents selectively choose what information to share.

3 Preliminaries

Coorperative multi-agent reinforcement learning. We formalize multi-agent reinforcement learning with a decentralized partially observable Markov decision process (Dec-POMDP), denoted as $\langle \mathcal{N}, \mathcal{S}, \{\mathcal{A}^i\}_{i\in\mathcal{N}}, P, R, \{O^i\}_{i\in\mathcal{N}}, \gamma\rangle$, where $\mathcal{N}:=\{1,\ldots,N\}$ represents the set of agents, \mathcal{S} the state space, \mathcal{A}^i the action space of agent i. At each time step t, each agent i selects an action $a_t^i \in \mathcal{A}^i$ to form a joint action $a_t = (a_1, \cdots, a_N) \in \prod_{i=1}^n \mathcal{A}^i =: \mathcal{A}$. The next state follows the distribution given by the state transition function $P: \mathcal{S} \times \mathcal{A} \to \Delta(\mathcal{S})$ as $s_{t+1} \sim P(\cdot \mid s_t, a_t)$, where $\Delta(\mathcal{X})$ is the set of probability distributions over set \mathcal{X} . All agents receive the same reward according to the reward function $R: \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ as $r_t := R(s_t, a_t)$. The observation function $O^i: \mathcal{S} \to \mathcal{O}^i$ generates an observation for agent i from its observation space \mathcal{O}^i , denoted as $o_t^i := O^i(s_t)$. Each agent chooses actions by sampling from its (fully decentralized) policy $\pi^i(a_t^i|\tau_t^i)$ that is conditioned on its action-observation trajectory $\tau_t^i := (o_0^i, a_0^i, o_1^i, a_1^i, ..., o_t^i)$. Agents' individual policies form the joint policy, $\pi:=(\pi^1,...,\pi^N)$, and their goal is to find π that maximizes expected discounted cumulative rewards $\mathbb{E}_\pi[\sum_{r=0}^\infty \gamma^t r_t]$.

Latent-space world models. Many recent works on world models rely on recurrent state-space models (RSSMs). We here review the core components of DreamerV2's RSSMs in the single-agent case. Initializing the latent state as h_0 , an RSSM encodes the agent's action-observation trajectory $(o_0, a_0, ..., o_{t-1}, a_{t-1})$ into the latent state h_t with a recurrent model f as $h_t = f(h_{t-1}, z_{t-1}, a_{t-1})$ where $z_t \sim q(z_t|h_t, o_t)$ is the encoded observation o_t (conditioned on h_t) by a (stochastic) representation model g. Paired with a (stochastic) transition model g are parameterized by neural networks and trained by maximizing the evidence lower bound (ELBO) of the log probability of $\log p(o_{0:T}|a_{0:T-1})$.

4 Method

4.1 Information sharing without bandwidth constraints

In the decentralized training and decentralized execution framework, each agent operates independently and accumulates distinct experiences. However, the constantly changing policies of other agents in the environment lead to inherent instability in the agent's learning process. To address this, we enable agents to exchange valuable information, so as to facilitate the optimization of their objective functions. Our approach first extends standard DreamerV2 to the MARL setting, then focuses on efficient information sharing and model optimization.

Centralized training. We consider a centralized training and decentralized execution framework as a performance reference to explore the potential of selective communication in fully decentralized MARL. Specifically, we adopt a shared world model and actor-critic network, parameterized by ϕ and θ , which are trained centrally using a replay buffer that aggregates information from all agents. Unlike centralized control, where a single policy outputs joint actions for all agents, execution in our approach remains decentralized, with each agent acting independently using the shared model. A similar parameter-sharing strategy has shown effectiveness in MABL (Venugopal et al., 2023), improving sample efficiency and overall performance in multi-agent tasks. This approach is based on DreamerV2 by incorporating RSSM, reconstruction models, and prediction models based on latent variables.

$$\text{RSSM:} \begin{cases} \text{Recurrent model:} & h_t^i = f_\phi(h_t^i \mid h_{t-1}^i, z_{t-1}^i, a_{t-1}^i) \\ \text{Representation model:} & z_t^i = q_\phi(z_t^i \mid o_t^i, h_t^i) \\ \text{Transition model:} & \hat{z}_t^i = p_\phi(\hat{z}_t^i \mid h_t^i) \end{cases}$$

The RSSM is a framework for learning latent dynamics. It consists of a recurrent model that maintains historical dependencies, a representation model that infers posterior distributions, and a transition model that predicts future states, facilitating synthetic trajectory generation.

In addition, we categorize the auxiliary components in the world model into two types: decoders and predictors. The decoders are responsible for decoding the latent representation back into actual behavior and perception, ensuring that the latent space contains sufficient environmental and agentspecific information. In contrast, predictors are used to infer future dynamics based on current latent states.

$$\begin{aligned} \text{Decoders:} & \begin{cases} \text{Observation decoder:} & \hat{o}_t^i = p_\phi(\hat{o}_t^i \mid h_t^i, z_t^i) \\ \text{Action decoder:} & \hat{a}_t^i = p_\phi(\hat{a}_t^i \mid h_t^i, z_t^i) \end{cases} \\ & \text{Reward predictor:} & \hat{r}_t^i = p_\phi(\hat{r}_t^i \mid h_t^i, z_t^i) \end{cases} \\ & \text{Predictors:} & \begin{cases} \text{Reward predictor:} & \hat{r}_t^i = p_\phi(\hat{r}_t^i \mid h_t^i, z_t^i) \\ \text{Termination predictor:} & \hat{r}_t^i = p_\phi(\hat{r}_t^i \mid h_t^i, z_t^i) \end{cases} \\ & \text{Available action predictor:} & \hat{A}_t^i = p_\phi(\hat{A}_t^i \mid h_t^i, z_t^i) \end{cases}$$

To generate more effective synthetic trajectories for policy optimization, predictors assist training by providing important feedback signals, including a reward predictor that outputs a continuous reward value \hat{r}_t^i , while the termination predictor outputs a binary variable $\hat{\gamma}_t^i$ to indicate whether the current state is terminal. The available action predictor outputs a vector \hat{A}_t^i of size A, where each element indicates whether the corresponding action is available at the time step t. Thus, both the termination and the available action predictors use Bernoulli distributions.

As described in the previous section, we optimize the decoders and predictors using supervised learning, and optimize the RSSM by maximizing the ELBO. For a trajectory of length T from agent i, the loss $L_{ELBO} = L_{\hat{o}t} + D_{KL}$ is computed as the expectation with respect to the posterior distribution $q_{\phi}(z_{1:T}^{i} \mid o_{1:T}^{i}, a_{1:T}^{i})$, to maximize the reconstruction accuracy and align the prior distribution p_{ϕ} with the posterior distribution q_{ϕ} . The loss functions are:

$$L_{\hat{o}_{t}} = -\sum_{i} \sum_{t} \log p_{\phi} (\hat{o}_{t}^{i} \mid h_{t}^{i}, z_{t}^{i}), \quad L_{\hat{a}_{t}} = -\sum_{i} \sum_{t} \log p_{\phi} (\hat{a}_{t}^{i} \mid h_{t}^{i}, z_{t}^{i})$$
(1)

$$L_{\hat{r}_{t}} = -\sum_{i} \sum_{t} \log p_{\phi} (\hat{r}_{t}^{i} \mid h_{t}^{i}, z_{t}^{i}), \quad L_{\hat{\gamma}_{t}} = -\sum_{i} \sum_{t} \log p_{\phi} (\hat{\gamma}_{t}^{i} \mid h_{t}^{i}, z_{t}^{i})$$
(2)

$$L_{\hat{A}_{t}} = -\sum_{i} \sum_{t} \log p_{\phi} (\hat{A}_{t}^{i} \mid h_{t}^{i}, z_{t}^{i}), \quad D_{KL} = \sum_{i} \sum_{t} KL[q_{\phi}(z_{t}^{i} \mid o_{t}^{i}, h_{t}^{i}) \parallel p_{\phi} (\hat{z}_{t}^{i} \mid h_{t}^{i})].$$
(3)

$$L_{\hat{r}_t} = -\sum_i \sum_t \log p_\phi(\hat{r}_t^i \mid h_t^i, z_t^i), \quad L_{\hat{\gamma}_t} = -\sum_i \sum_t \log p_\phi(\hat{\gamma}_t^i \mid h_t^i, z_t^i)$$
 (2)

$$L_{\hat{A}_t} = -\sum_{i} \sum_{t} \log p_{\phi}(\hat{A}_t^i \mid h_t^i, z_t^i), \quad D_{KL} = \sum_{i} \sum_{t} KL[q_{\phi}(z_t^i \mid o_t^i, h_t^i) \parallel p_{\phi}(\hat{z}_t^i \mid h_t^i)]. \quad (3)$$

The model-based approach effectively decouples model learning from policy learning in MARL. Once the world model is trained, it generates imagined trajectories for policy optimization. We employ an actor-critic framework to enhance agents' decision-making and coordination. Each agent utilizes the shared parameter policy network π_{θ} , and the objective is to maximize the cumulative MARL returns, thereby learning an optimal policy. Specifically, at time step t, the agent selects an action based on the following:

$$a_t^i \sim \pi_\theta(a_t^i \mid \hat{z}_t^i, h_t^i). \tag{4}$$

Here, the agent performs policy inference based on its own hidden state vector h_t^i and the inferred prior distribution \hat{z}_t^i . And the shared critic V_ϕ estimates each agent's value function to guide policy optimization based on:

$$\hat{V}_t^i \sim V_\phi(\hat{z}_t^i, h_t^i). \tag{5}$$

In this centralized training method, all agents share the parameters of the world models, actor and critic networks, enabling efficient learning and coordination for multi-agent systems.

Independent training. In order to verify the impact of information sharing in fully decentralized MARL, we introduce a lower bound method, which follows an independent training paradigm. Initially, all agents have identical architectures and optimizers for the world models and actor-critic models, along with the same parameters ϕ^i_0 and θ^i_0 . During training, each agent treats other agents as part of their environment, builds its own world model, optimizes policy without sharing parameters or information, and follows the previously introduced loss functions but iterating solely over time steps.

Sharing experiences across the RSSM and predictors. We first propose the RSSM+Predictors method, in which part of the world models share experiences while others remain independent. Specifically, RSSM and predictors are trained by sampling from a shared experience replay buffer used by all agents. This global perspective allows each agent not only to rely on its own experience for training but also to capture global dynamics and environmental features across agents, thereby enhancing the understanding of environmental patterns. On the other hand, each agent's decoders depend entirely on its own independent experience, ensuring that each agent can optimize its decoders based on its own perspective and behavioral patterns, thus preventing a decline in decoding accuracy due to differences in experiences. Through this design, the RSSM and predictors focus on learning general latent-space representations and predicting accurate environmental dynamics, while the decoders focus on generating predictions of observations and actions based on the agent's own experience.

Algorithm 1 illustrates how experiences are shared for training the RSSM and predictor models. The method contains two types of buffers: one shared experience buffer and N independent experience buffers. During the experience collection phase, each agent stores its observations, actions, rewards, and other feedback information into the shared buffer and its independent buffer. During training, the RSSM and predictors update using samples from the shared buffer, while the decoders optimize only with the agent-specific buffer.

Algorithm 1 Training World Model with Shared and Individual Experience Buffers

```
1: Initialize SharedReplayBuffer
 2: Initialize N IndividualReplayBuffers
 3: for t = 1 to T do
          for each agent i = 1 to N do
 4:
               SharedReplayBuffer.Add((o_t^i, a_t^i, r_t^i, \gamma_t^i, A_t^i))
 5:
               IndividualReplayBuffers[i].Add((o_t^i, a_t^i, r_t^i, \gamma_t^i, A_t^i))
 6:
 7:
          end for
          for each agent i = 1 to N do
 8:
                GradientStep(\phi^i, Equation 1) on samples from IndividualReplayBuffers[i]
 9:
                GradientStep(\phi^i, Equation 2 to 3) on samples from SharedReplayBuffer
10:
                IndividualReplayBuffers[i].Sample(batch_size)
11:
               \begin{array}{l} \hat{a}_{1:H}^i, \hat{A}_{1:H}^i, logit_{1:H}^i, \hat{r}_{1:H}^i, h_{1:H}^i, z_1^i, \hat{z}_{2:H}^i = \text{ImaginationRollout}(o_t^i, a_t^i, \gamma_t^i) \\ \text{GradientStep}[(\theta^i, \text{Equation 4}), (\phi^i, \text{Equation 5})] \end{array}
12:
13:
          end for
14:
15: end for
```

Sharing rollouts across the actor-critic networks. In this method, the optimization of actor and critic networks is achieved through rollouts sharing. Specifically, each agent trains its independent world model based on its own experience and generates synthetic trajectories over a certain horizon with length H. As shown in Algorithm 2, these trajectories include not only sequences of actions $\hat{a}_{1:H}^i$ and rewards $\hat{r}_{1:H}^i$, but also incorporate latent state information $h_{1:H}^i$, z_1^i , and $\hat{z}_{2:H}^i$, and auxiliary information $\hat{A}_{1:H}^i$ and $logit_{1:H}^i$, where $logit_h^i$ denotes the logits corresponding to \hat{a}_h^i , ensuring dif-

ferentiability. Subsequently, all agents share their generated synthetic trajectories and the aggregated multi-source trajectory data is used for training their actor-critic networks.

On one hand, fully utilizing model-generated data alleviates the issue of high interaction costs, improving sample efficiency and exploration capability. On the other hand, this sharing mechanism enables agents to perform policy learning in a wider data distribution, thereby enhancing decision-making and accelerating the convergence of the training process.

Algorithm 2 Training Actor-Critic with Shared Imagination Buffers

```
1: Initialize N IndividualReplayBuffers
 2: Initialize SharedImaginationBuffer
 3: for t = 1 to T do
            for each agent i = 1 to N do
 4:
 5:
                  IndividualReplayBuffers[i].Sample(batch_size)
                   GradientStep(\phi^i, Equation 1 to 3)
 6:
                  \begin{split} \hat{a}_{1:H}^{i}, \hat{A}_{1:H}^{i}, logit_{1:H}^{i}, \hat{r}_{1:H}^{i}, h_{1:H}^{i}, z_{1}^{i}, \hat{z}_{2:H}^{i} &= \text{ImaginationRollout}(o_{t}^{i}, a_{t}^{i}, \gamma_{t}^{i}) \\ \text{SharedImaginationBuffer.Add}(\hat{a}_{1:H}^{i}, \hat{A}_{1:H}^{i}, logit_{1:H}^{i}, \hat{r}_{1:H}^{i}, h_{1:H}^{i}, z_{1}^{i}, \hat{z}_{2:H}^{i}) \end{split}
 7:
 8:
            end for
 9:
            for each agent i = 1 to N do
10:
                  SharedImaginationBuffer.Sample(batch_size)
11:
                  GradientStep[(\theta^i, \text{Equation 4}), (\phi^i, \text{Equation 5})]
12:
            end for
13:
14: end for
```

4.2 Information sharing with bandwidth constraints

Although MARL agents can communicate without restriction, unconstrained communication may introduce excessive redundancy, leading to communication overload and unstable parameter updates during training. Gerstgrasser et al. (2023) have shown that selectively sharing experiences, rather than indiscriminately broadcasting them, can significantly accelerate the learning process. Motivated by these observations, we introduce a bandwidth-aware information sharing strategy to improve communication efficiency and training stability. Specifically, whenever new experiences or imaginations are collected, each agent not only updates its own replay buffer, but also selects the most relevant information for sharing, which are subsequently integrated into the replay buffers of other agents. In the previous section, we introduced the idea of selectively sharing information from the RSSM, predictors, and actor-critic components. In this section, we further consider adaptive control of information sharing under explicit bandwidth constraints, thus improving data sharing efficiency and optimizing communication resource utilization.

During world model training, experience sharing involves four key component models. However, due to significant variations in the loss distribution and scale across different models, applying a unified standard directly would lead to imbalanced sharing criteria. Therefore, we first use a multimodel loss normalization and aggregation strategy, ensuring fair and stable sample selection. At each training step, we first compute the loss values for the four models: $L_{\hat{A}_t}$, $L_{\hat{\gamma}_t}$, $L_{\hat{\tau}_t}$, D_{KL} . Then we maintain historical statistics for each loss function within a sliding window of length K, including the mean and standard deviation:

$$\mu_m = \frac{1}{K} \sum_{k=0}^{K} L_k, \qquad \sigma_m = \sqrt{\frac{1}{K} \sum_{k=0}^{K} (L_k - \mu_m)^2}.$$

Each loss is normalized as $L_m' = \frac{L_m - \mu_m}{\sigma_m + \epsilon}$, where ϵ is a small constant to prevent division by zero.

After normalization, we aggregate the four loss values into a single composite loss L_M for sample selection: $L_M = \max(L'_{\hat{A}_t}, L'_{\hat{\gamma}_t}, L'_{\hat{\gamma}_t}, D'_{KL})$. This strategy ensures that if any single model exhibits

an abnormally high loss, the corresponding sample is selected for sharing, reducing the risk of selection being dominated by a single model's loss and ensuring balanced multi-model learning. To adaptively control the number of shared samples, we adopt the deterministic Gaussian experience selection from Gerstgrasser et al. (2023) over the quantile-based alternative due to its efficiency and robustness. It enables bandwidth control through a tunable parameter β without requiring sorting or quantile computation. Additionally, it only requires tracking the running mean and variance, without storing the full cluster of recent samples. Specifically, we maintain the most recent K samples to track the distribution of the composite loss L_M , including its mean μ_M and standard deviation σ_M . Then an experience is selected for sharing if:

$$L_M > \mu_M + c \cdot \sigma_M$$

where c is a constant determined by the target bandwidth β , satisfying $1 - \text{cdf}N(c) = \beta$. This enables agents to adaptively select the top β -fraction of their most informative experiences.

We extend this selection strategy to the actor-critic setting, focusing solely on the critic loss as the primary criterion for sample selection. Experimental results (Figure 3) show that critic loss exhibits significant variation between different methods (Actor-Critic, Centralized Training, Independent Training), indicating that training strategies mainly impact the critic network. Therefore, critic loss is a reliable indicator of agent learning progress and policy evaluation accuracy. In contrast, actor loss exhibits less variation in both magnitude and trend in different methods. Based on these observations, we adopt a critic loss based selection method in the actor-critic training process. During training, we apply a similar strategy as used for the world model but exclusively utilize critic loss for selecting highly informative samples. The statistics of the critic loss in the sliding window are calculated as:

$$\mu_{\text{critic}} = \frac{1}{K} \sum_{k=0}^{K} L_k, \qquad \sigma_{\text{critic}} = \sqrt{\frac{1}{K} \sum_{k=0}^{K} (L_k - \mu_{\text{critic}})^2}.$$

Then, it applies imagined rollouts sharing based on the following inequality:

$$L_{\text{critic}} \ge \mu_{\text{critic}} + c \cdot \sigma_{\text{critic}}.$$

5 Experiments

To evaluate decentralized cooperation in multi-agent systems with experience and imagination rollout sharing, we use the StarCraft Multi-Agent Challenge (SMAC) benchmark, based on StarCraft II (Vinyals et al., 2017; Samvelyan et al., 2019). Each agent independently executes its policy while coordinating with other agents to defeat enemy units. Specifically, we conduct experiments on two micro-trick scenarios in which two agents face a single enemy (2s_vs_1sc and 2m_vs_1z), a homogeneous and symmetric scenario where both armies consist of three Marines (3m), and a heterogeneous and symmetric scenario where the allied team consists of two Stalkers and three Zealots, matching the enemy composition (2s3z). In each scenario, we perform ten independent runs for each method with an equal number of training steps. To ensure fair comparison, we ensure that the architecture and setup of all models are identical across all methods.

Baselines: In our experiments, we compare the proposed method with the Independent Training baseline to assess the impact of information sharing in decentralized cooperation and also to find out how it compares to the Centralized Training baseline.

The Independent Training baseline represents a fully decentralized variant of DreamerV2, where each agent is trained independently without any parameter sharing or inter-agent communication. At the other extreme, the Centralized Training baseline serves as an upper bound on information sharing: all agents' experiences are aggregated into a shared replay buffer and used to jointly train a single world model and actor-critic network.

In contrast, our method follows a decentralized training and decentralized execution framework. Each agent maintains and updates its own model independently while selectively sharing a small subset of high-value experiences or imagined rollouts with others. These correspond to the methods RSSM+Predictors and Actor-Critic, respectively. This comparison allows us to systematically analyze the benefits of bandwidth-constrained partial communication under decentralized learning dynamics.

Implementation details: Our work builds on the official implementation of MAMBA. All neural modules use fully connected feedforward networks with ReLU activations, including the world model, actor, and critic. We adopt standard hyperparameter settings with moderate adjustments tailored to our experimental setup. Specifically, the sequence length is set to 20, and the rollout horizon is set to 15 steps. The replay buffer stores up to 250,000 transitions, and a sliding window of size 1500 is used for information selection, enabling Gaussian-based filtering over recent samples for bandwidth-constrained sharing. Each iteration consists of five epochs for the world model and four epochs for the actor-critic updates, with a batch size of 40. The entropy coefficient is set to 0.01 and decays with a factor of 0.99998. We apply gradient clipping with a norm of 100 for the world model and 20 for the policy network. All hidden layers are set to a dimensionality of 256. The learning rates are configured as 1e-6 for the actor, 1e-4 for the critic, and 3e-4 for the world model.

In the following, we investigate the impact of selective information sharing in decentralized MARL. Section 5.1 evaluates the effectiveness of selectively sharing information across the RSSM, predictors, and actor-critic components under decentralized training, compared to the two baselines. Section 5.2 explores how varying the communication bandwidth affects the learning process by controlling the amount of information shared among agents. We provide a detailed analysis of communication efficiency and training performance under the proposed selective sharing strategies.

5.1 Performance comparison without bandwidth constraints

To evaluate the effectiveness of two information sharing strategies, we first compare the performance of different methods in four SMAC scenarios without considering bandwidth constraints. Figure 1 presents the win rate curves of these four methods over training steps in different scenarios. To further analyze the underlying factors contributing to performance improvements, Figure 2 and Figure 3 present the loss curves for RSSM+Predictors and Actor-Critic models of different methods, respectively, in the 2s_vs_1sc scenario. And we choose to discuss only the loss curves of Agent-0 for the RSSM+Predictors models, as all agents exhibit similar patterns in the allied team.

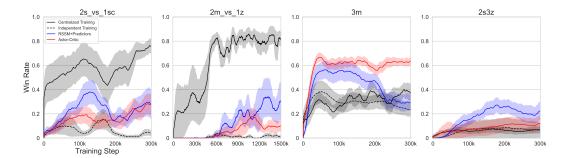


Figure 1: Win rate curves.

As shown in Figure 1, we can see that RSSM+Predictors (experience sharing) and Actor-Critic (imagined rollouts sharing) methods consistently outperform Independent Training in all scenarios, achieving higher win rates and faster convergence. This improvement highlights the effectiveness of the loss-aware information sharing mechanism in facilitating superior coordination strategies, whereas Independent Training struggles to achieve high win rates under fully decentralized learning. In the micro-trick asymmetric scenarios (2s_vs_1sc and 2m_vs_1z), agents need precise

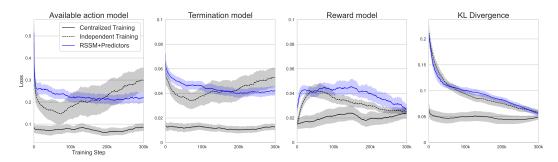


Figure 2: Loss curves of an agent on 2s_vs_1sc of different methods for RSSM+Predictors models.

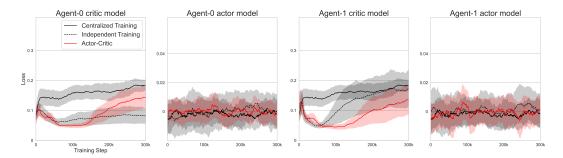


Figure 3: Loss curves on SMAC of different methods for Actor-Critic models.

coordination to attack fewer enemies than themselves. Centralized Training enables all agents to learn and perform best through parameter sharing, thus showing excellent performance. However, in RSSM+Predictors and Actor-Critic methods, agents only share information and still optimize strategies independently, which makes it difficult to fully match the overall collaboration ability of Centralized Training, resulting in suboptimal performance. Independent Training, due to its complete decentralization, agents cannot learn to cooperate effectively, resulting in the lowest win rate. For the homogeneous and heterogeneous symmetric scenarios (3m and 2s3z), we conjecture that Centralized Training may lead to unnecessary synchronous behavior due to parameter sharing. This could result in all agents attacking the same target simultaneously or retreating at the same time, potentially reducing combat efficiency. In contrast, RSSM+Predictors and Actor-Critics methods allow agents to share information while independently optimizing their policies; this might make them more adaptable when executed in a decentralized manner, thus surpassing Centralized Training. As shown in Figure 2, in 2s vs 1sc scenario, the loss curves for the available action model, termination model, reward model and the KL divergence curves maintain lower loss values throughout training compared to Independent Training, even approaching the loss levels of Centralized Training. This indicates that experience sharing stabilizes training and improves prediction accuracy. Specifically, the KL divergence results show that RSSM+Predictors achieves better latent space alignment, demonstrating the effectiveness of experience sharing in approximating centralized training.

The actor-critic loss curves shown in Figure 3 provide another insight into the training dynamics of our method. For 2s_vs_1sc scenario, the critic loss curves for both Agent-0 and Agent-1 show distinct trends, whereas the actor loss curves remain similar in magnitude and trend across all methods. Specifically, the Independent Training method maintains a relatively lower critic loss while the Centralized Training method exhibits a rising critic loss over time, indicating potential overfitting or instability in value estimation due to indiscriminate sharing of imagined rollouts. And our Actor-Critic method shows an intermediate performance between the Centralized Training and Independent Training methods that aligns with the win rate trends shown in Figure 1. For the actor, the loss curves show negligible differences between methods and maintain relatively stable training trends. This suggests that behavior learning is less affected by the rollouts sharing compared to value function estimation.

Overall, these results indicate that both experience and imagined rollouts sharing can benefit the multi-agent coordination in decentralized training. Performance gains can be attributed to the use of training loss as a criterion for information sharing, enabling agents to generalize and coordinate effectively without relying on shared parameters.

5.2 Performance comparison with bandwidth constraints

In this section, we further investigate the impact of information sharing under limited bandwidth by evaluating RSSM+Predictors and Actor-Critic across varying target bandwidths in the 2s_vs_1sc and 3m scenarios. As shown in Figure 4, the Centralized Training method achieves the highest performance, as it benefits from parameter sharing. In comparison, the Independent Training method performs as a lower bound, reflecting the challenges of fully decentralized learning without any inter-agent communication. In particular, considering the limited bandwidth, we can observe a clear peak in RSSM+Predictors and Actor-Critic methods around the target bandwidth of 0.01 and 0.05 respectively for 2s_vs_1sc scenario. The results indicate that selective information sharing can significantly promote the learning process compared to Independent Training. Interestingly, the ShareAll setting, where all information is shared without selection, does not necessarily lead to better performance compared to intermediate bandwidth values. This suggests that excessive communication may introduce redundant information, potentially hindering learning efficiency. In contrast, our Gaussian-based selection mechanism prioritizes information based on their prediction and estimation quality, enabling more efficient and focused communication.

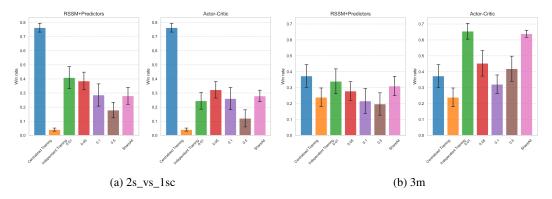


Figure 4: Comparison of RSSM+Predictors and Actor-Critic with different bandwidths.

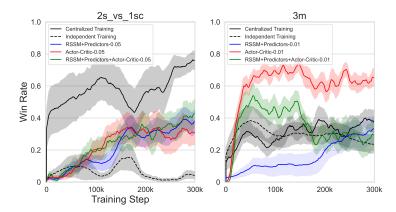


Figure 5: Comparison of RSSM+Predictors and Actor-Critic with optimal target bandwidths.

Figure 5 further compares the methods in their best-performing bandwidths. Specifically, based on Figure 4a and Figure 4b, we adopt the optimal target bandwidth of 0.05 and 0.01 for the 2s_vs_1sc

and 3m scenarios, respectively. In 2s_vs_1sc, both RSSM+Predictors-0.05 and Actor-Critic-0.05 significantly outperform Independent Training, demonstrating the effectiveness of decentralized information sharing. However, their combination does not lead to a significant additional gain, with performance remaining comparable to that of each method individually. This suggests that the two approaches may provide overlapping rather than complementary benefits. RSSM+Predictors contributes compact latent representations that enhance predictive accuracy, while Actor-Critic shares imagined rollouts that support strategic foresight. In this asymmetric environment, either type of information alone may satisfy the coordination requirements sufficiently. In contrast, in 3m, a symmetric environment where agents perform identical roles and face mirrored opponents, both RSSM+Predictors-0.01 and Actor-Critic-0.01 outperform Independent Training. In particular, Actor-Critic-0.01 achieves the highest win rate, likely because shared imagined rollouts provide greater benefits in tasks requiring precise behavioral alignment among agents. RSSM+Predictors-0.01, while still helpful, is less impactful when agents possess similar observations and responsibilities.

These findings indicate that the effectiveness of different information sharing strategies is highly dependent on the coordination structure of the environment. Our Gaussian-based selection mechanism facilitates this adaptation by filtering shared information based on prediction or estimation quality, allowing each method to focus on transmitting only the most relevant information aligned with its strengths.

6 Conclusion

In this paper, we investigate information sharing strategies in model-based MARL under a decentralized training and decentralized execution framework. We propose selective communication mechanisms for both the world model and the actor-critic training processes. Specifically, the world model is decomposed so that the RSSM and prediction modules benefit from shared experiences, while the decoders remain agent-specific to preserve individual features. We then evaluate the performance of the resulting RSSM+Predictors and Actor-Critic methods across varying bandwidth constraints.

Our experimental results in four SMAC scenarios show that both methods significantly outperform Independent Training across a range of bandwidths and, in some cases, even surpass the performance of Centralized Training, demonstrating the effectiveness of structured information sharing. Interestingly, we find that intermediate bandwidth levels often lead to the best performance, suggesting that excessive communication can introduce redundancy or noise and is not always beneficial. Moreover, the utility of shared information is influenced by environment symmetry. In asymmetric tasks, either experiences or imagined rollouts can enhance coordination, while in symmetric tasks, shared imagined rollouts are more effective in synchronizing agent behavior. These findings also indicate that combining RSSM+Predictors and Actor-Critic at their respective optimal bandwidths does not yield additive improvements, but rather reflects overlapping contributions. Our Gaussian-based selection mechanism further enhances both methods by prioritizing the most relevant information, enabling more targeted and scalable communication.

Acknowledgments

The authors acknowledge funding support from NSF award IIS-2154904. The authors thank the anonymous reviewers for their insightful and constructive reviews.

References

Dingyang Chen, Yile Li, and Qi Zhang. Communication-efficient actor-critic methods for homogeneous markov games. *arXiv preprint arXiv:2202.09422*, 2022.

Filippos Christianos, Lukas Schäfer, and Stefano Albrecht. Shared experience actor-critic for multiagent reinforcement learning. *Advances in neural information processing systems*, 33:10707–10717, 2020.

- Tianshu Chu, Jie Wang, Lara Codecà, and Zhaojian Li. Multi-agent deep reinforcement learning for large-scale traffic signal control. *IEEE transactions on intelligent transportation systems*, 21(3): 1086–1095, 2019.
- Peter Corke, Ron Peterson, and Daniela Rus. Networked robots: Flying robot navigation using a sensor net. In *Robotics research. The eleventh international symposium*, pp. 234–243. Springer, 2005.
- Joris Dinneweth, Abderrahmane Boubezoul, René Mandiau, and Stéphane Espié. Multi-agent reinforcement learning for autonomous vehicles: A survey. *Autonomous Intelligent Systems*, 2(1):27, 2022.
- Vladimir Egorov and Aleksei Shpilman. Scalable multi-agent model-based reinforcement learning. arXiv preprint arXiv:2205.15023, 2022.
- Matthias Gerstgrasser, Tom Danino, and Sarah Keren. Selectively sharing experiences improves multi-agent reinforcement learning. *Advances in Neural Information Processing Systems*, 36: 59543–59565, 2023.
- Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning behaviors by latent imagination. *arXiv preprint arXiv:1912.01603*, 2019a.
- Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and James Davidson. Learning latent dynamics for planning from pixels. In *International conference on machine learning*, pp. 2555–2565. PMLR, 2019b.
- Danijar Hafner, Timothy Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering atari with discrete world models. *arXiv preprint arXiv:2010.02193*, 2020.
- Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains through world models. *arXiv preprint arXiv:2301.04104*, 2023.
- Pablo Hernandez-Leal, Michael Kaisers, Tim Baarslag, and Enrique Munoz De Cote. A survey of learning in multiagent environments: Dealing with non-stationarity. *arXiv* preprint *arXiv*:1707.09183, 2017.
- Woojun Kim, Jongeui Park, and Youngchul Sung. Communication in multi-agent reinforcement learning: Intention sharing. In *International conference on learning representations*, 2020.
- Ryan Lowe, Yi I Wu, Aviv Tamar, Jean Harb, OpenAI Pieter Abbeel, and Igor Mordatch. Multiagent actor-critic for mixed cooperative-competitive environments. *Advances in neural information processing systems*, 30, 2017.
- Frans A Oliehoek, Christopher Amato, et al. *A concise introduction to decentralized POMDPs*, volume 1. Springer, 2016.
- Tabish Rashid, Mikayel Samvelyan, Christian Schroeder De Witt, Gregory Farquhar, Jakob Foerster, and Shimon Whiteson. Monotonic value function factorisation for deep multi-agent reinforcement learning. *Journal of Machine Learning Research*, 21(178):1–51, 2020.
- Mikayel Samvelyan, Tabish Rashid, Christian Schroeder De Witt, Gregory Farquhar, Nantas Nardelli, Tim GJ Rudner, Chia-Man Hung, Philip HS Torr, Jakob Foerster, and Shimon Whiteson. The starcraft multi-agent challenge. *arXiv preprint arXiv:1902.04043*, 2019.
- Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Simon Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, et al. Mastering atari, go, chess and shogi by planning with a learned model. *Nature*, 588(7839):604–609, 2020.
- Richard S Sutton. Dyna, an integrated architecture for learning, planning, and reacting. *ACM Sigart Bulletin*, 2(4):160–163, 1991.

- Edan Toledo and Amanda Prorok. Codreamer: Communication-based decentralised world models. *arXiv preprint arXiv:2406.13600*, 2024.
- Aravind Venugopal, Stephanie Milani, Fei Fang, and Balaraman Ravindran. Mabl: Bi-level latent-variable world model for sample-efficient multi-agent reinforcement learning. *arXiv* preprint arXiv:2304.06011, 2023.
- Oriol Vinyals, Timo Ewalds, Sergey Bartunov, Petko Georgiev, Alexander Sasha Vezhnevets, Michelle Yeo, Alireza Makhzani, Heinrich Küttler, John Agapiou, Julian Schrittwieser, et al. Starcraft ii: A new challenge for reinforcement learning. *arXiv preprint arXiv:1708.04782*, 2017.
- Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David H Choi, Richard Powell, Timo Ewalds, Petko Georgiev, et al. Grandmaster level in starcraft ii using multi-agent reinforcement learning. *nature*, 575(7782):350–354, 2019.
- Zhiwei Xu, Bin Zhang, Yuan Zhan, Yunpeng Baiia, Guoliang Fan, et al. Mingling foresight with imagination: Model-based cooperative multi-agent reinforcement learning. *Advances in Neural Information Processing Systems*, 35:11327–11340, 2022.
- Weirui Ye, Shaohuai Liu, Thanard Kurutach, Pieter Abbeel, and Yang Gao. Mastering atari games with limited data. *Advances in neural information processing systems*, 34:25476–25488, 2021.