

Hierarchical Message-Passing Policies for Multi-Agent Reinforcement Learning

Tommaso Marzi¹, Cesare Alippi^{1,2}, Andrea Cini^{3,1}

¹ IDSIA USI-SUPSI, Università della Svizzera italiana, Lugano, Switzerland

² Politecnico di Milano, Milan, Italy

³ Swiss National Science Foundation Postdoctoral Fellow

{tommaso.marzi, cesare.alippi, andrea.cini}@usi.ch

Abstract

Decentralized Multi-Agent Reinforcement Learning (MARL) methods allow for learning scalable multi-agent policies, but suffer from partial observability and induced non-stationarity. These challenges can be addressed by introducing mechanisms that facilitate coordination and high-level planning. Specifically, coordination and temporal abstraction can be achieved through communication (e.g., message passing) and Hierarchical Reinforcement Learning (HRL) approaches to decision-making. However, optimization issues limit the applicability of hierarchical policies to multi-agent systems. As such, the combination of these approaches has not been fully explored. To fill this void, we propose a novel and effective methodology for learning multi-agent hierarchies of message-passing policies. We adopt the feudal HRL framework and rely on a hierarchical graph structure for planning and coordination among agents. Agents at lower levels in the hierarchy receive goals from the upper levels and exchange messages with neighboring agents at the same level. To learn hierarchical multi-agent policies, we design a novel reward-assignment method based on training the lower-level policies to maximize the advantage function associated with the upper levels. Results on relevant benchmarks show that our method performs favorably compared to the state of the art.

1 Introduction

Designing information processing methods to support coordination and planning is a core challenge in Multi-Agent Reinforcement Learning (MARL) [48, 33]. Purely centralized approaches [29] reduce the multi-agent problem to a single-agent formulation and thus cannot scale [74]. Conversely, decentralized methods [6] rely only on local observations. While decentralized approaches are scalable, operating at the level of individual agents can make the environment appear non-stationary as other agents update their policies [31, 58]. The induced non-stationarity, combined with partial observability, can hinder coordination and high-level planning [36]. Similarly, Centralized Training Decentralized Execution (CTDE) approaches [24, 74, 6] train agents on all the available (global) information, but are limited to relying exclusively on local observations at execution time [75, 50]. To overcome these challenges, several methods allow for information sharing among agents. The Communication Multi-Agent Deep RL (Comm-MADRL) framework [76] improves coordination and mitigates partial observability by leveraging agent communication [65, 17, 32], possibly during both training and execution [57, 56]. A widely used approach to implement communication in MARL is to combine it with graph-based representations [11], which model communicating agents as *connected* nodes within a graph structure [20, 21]. In graph-based MARL, local observations are usually processed by assuming homogeneous agents and relying on weight-sharing message-

passing Graph Neural Networks (GNNs) [26, 8, 16]. However, communication on *flat* graphs (i.e., graphs simply modeling pairwise interactions) can introduce bottlenecks in information propagation and planning [44, 54]. Conversely, Hierarchical Reinforcement Learning (HRL) [68, 53, 10, 25] methods organize the learning system into a hierarchical decision-making structure. By doing so, HRL introduces learning biases that facilitate spatiotemporal abstraction and long-term planning [43, 72]. In particular, Feudal Reinforcement Learning (FRL) [18] relies on hierarchies of policies where the upper levels send commands and rewards to lower levels. In FRL, the optimization is based on level-specific reward signals, with the highest level receiving the reward directly from the environment. Combining communication-based approaches with a HRL framework such as FRL has the potential for bringing together the best of both worlds, but is not straightforward. In particular, FRL relies on designing ad-hoc reward functions for the lower levels. This adds overhead for the designer and makes end-to-end learning of the hierarchy of policies challenging [69, 54]. Furthermore, adding a hierarchy on top of the base agents requires the adaptation of message-passing mechanisms for the additional interacting entities.

In this paper, we propose a novel HRL method for communication and coordination among different policies in MARL problems. We do so by adopting the FRL framework and introducing a procedure to avoid the associated reward-shaping and optimization issues. Our approach allows for learning a multi-level hierarchy of message-passing policies based on FRL [18, 54], thereby achieving high-level planning and coordination. The hierarchical graph defining the feudal structure is determined according to the current state, which can change over time. To train the resulting hierarchy of policies, we design a novel reward-assignment scheme based on training the lower-level policies to maximize the advantage function associated with the upper levels. The resulting propagation mechanism is decentralized as each agent optimizes a different and separate reward function dependent on the structure of the hierarchy at each time step. In particular, we theoretically show that the resulting objective remains aligned with maximizing agents’ local reward. This also makes our method suitable for scenarios that are not fully cooperative (provided a coherent definition of the hierarchical structure). Our hierarchical MARL approach enjoys the same benefits of decentralized communication methods [76, 5], while relying on the feudal paradigm to mitigate non-stationarity and partial observability. Policy optimization can be carried out by applying any policy gradient algorithm to local trajectories at different levels. To summarize, our *contributions* are the following:

1. We introduce a novel method based on FRL for learning multi-level hierarchies of message-passing policies in multi-agent systems.
2. We propose a flexible and adaptive reward-assignment scheme that leverages hierarchical graph structures for training multi-level feudal policies.
3. We provide theoretical guarantees that the proposed learning scheme generates level-specific reward signals aligned with the global task.

Empirical results show that our method, named *Hierarchical Message-Passing Policy* (HiMPo), achieves strong performance in challenging MARL benchmarks where coordination among agents is required. HiMPo does not necessarily require a pre-defined graph structure to operate on, as it can simply rely on the hierarchy for communication. HiMPo is a step toward the design of effective MARL methods for coordination and high-level planning.

2 Preliminaries

Partially-Observable Stochastic Game Multi-agent systems can be formalized as Partially-Observable Stochastic Games (POSGs) [4], i.e., as tuple $\langle \mathcal{N}, \mathcal{S}, \{\mathcal{Z}_i\}_{i \in \mathcal{N}}, \{\mathcal{A}_i\}_{i \in \mathcal{N}}, \mathcal{P}, \{\mathcal{R}_i\}_{i \in \mathcal{N}}, \gamma \rangle$ where \mathcal{N} is the set of agents, \mathcal{S} is a state space, \mathcal{Z}_i and \mathcal{A}_i are the observation and action spaces of the i -th agent, $\mathcal{P} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$ is a Markovian transition function depending on the joint set of actions, $\mathcal{R}_i : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the payoff (reward) function for the agent i , and $\gamma \in [0, 1)$ is a discount factor. Local observations $\{\mathbf{z}_t^i\}_{i \in \mathcal{N}}$ are generated as a function of the global state \mathbf{S}_t . We consider *homogeneous* POSGs [27], i.e., systems where agents have identical (but individual) observation spaces, action spaces, and reward functions. For each agent, the goal is to maximize the discounted sum of local rewards. If the task is fully cooperative, agents share the same objective. Conversely, in mixed (cooperative-competitive) settings, agents’ decisions are also driven by self-interested motives, necessitating local rewards [4].

Hierarchies of Feudal Policies Feudal Reinforcement Learning [18] approaches organize the decision-making structure into a multi-level hierarchy of policies. In FRL, actions taken by lower levels are conditioned on the goals assigned by upper levels, which are also responsible for propagating reward signals to subordinate agents. In the Feudal Graph Reinforcement Learning [54] paradigm, feudal dependencies are represented with a hierarchical graph \mathcal{G}^* , each node corresponding to a level-specific entity. In particular, (sub-)managers are abstract nodes sending goals to their subordinates, while workers represent sub-components of the agent (e.g., actuators). Each level maximizes a separate payoff signal: the top-level manager collects rewards directly from the environment, while lower levels receive an intrinsic reward based on the alignment with received commands.

3 Learning Multi-Agent Hierarchical Message-Passing Policies

Existing architectures for learning hierarchies of feudal policies require designing ad-hoc reward mechanisms, which can result in optimization issues [69, 54]. Applying such methodologies directly to multi-agent systems would suffer from the same drawbacks. Furthermore, in MARL, training the collection of policies by relying on a shared (global) reward neglects the individual contribution of each agent (besides assuming a fully cooperative scenario). To address these challenges, HiMPo learns a feudal hierarchy of message-passing policies in a decentralized fashion by independently optimizing level-specific signals. In HiMPo, the rewards for lower levels are defined based on the advantage function of the upper levels, computed using local returns. Different from previous works on multi-agent HRL [70, 49], the contribution of each action (goal) is evaluated locally w.r.t. each agent (manager) rather than globally. The following section introduces the proposed methodology. As reference, we consider a 3-level hierarchical graph (refer to Fig. 1, left). In particular, after providing an overview of the methodology, we illustrate how to build the hierarchical graph \mathcal{G}_t^* and use it to assign the level-specific rewards. Then, we provide details on the training routine and theoretical guarantees on the soundness of the proposed reward scheme.

3.1 Overview

At each step t , we represent the multi-level decision-making process through a dynamic hierarchical graph \mathcal{G}_t^* . We denote node-level representations at round l of message passing as $\mathbf{h}_t^{i,l}$, where i can refer to a worker ($i = w$), sub-manager ($i = s$), or top-level manager ($i = m$). Workers' (agents') representations $\{\mathbf{h}_t^{w,0}\}_{w \in \mathcal{N}}$ are initialized from the raw observations $\{\mathbf{z}_t^i\}_{i \in \mathcal{N}}$. These representations are then propagated bottom-up through \mathcal{G}_t^* to initialize the representation $\mathbf{h}_t^{s,0}$ of each sub-manager s . Finally, the representation $\mathbf{h}_t^{m,0}$ of the top-level manager m is initialized based on sub-managers' representations. Information is propagated by performing L_r message-passing rounds [26] w.r.t. the graph at each level of the hierarchy as

$$\mathbf{h}_t^{i,l+1} = \text{UPDATE}_l^i \left(\mathbf{h}_t^{i,l}, \text{AGGR}_{j \in \mathcal{B}_t(i)} \left\{ \text{MSG}_l^i(\mathbf{h}_t^{i,l}, \mathbf{h}_t^{j,l}) \right\} \right), \quad (1)$$

where $\mathcal{B}_t(i)$ denotes the neighbors of the i -th node at time step t , UPDATE_l^i and MSG_l^i indicate learnable update and message function (e.g., MLPs), and AGGR is a permutation-invariant aggregation function (e.g., the mean). In practice, update and message functions are shared among nodes belonging to the same level. Since the top-level manager m is a single node without neighbors at the highest level of the hierarchy, we set $\mathbf{h}_t^m \equiv \mathbf{h}_t^{m,0}$. Then, based on such representations, the top-level manager m and each s -th sub-manager send goals ($\mathbf{g}_t^{m \rightarrow s} \in \mathbb{R}^{d_m}$ and $\mathbf{g}_t^{s \rightarrow w} \in \mathbb{R}^{d_s}$, respectively), and each w -th worker (agent) takes actions ($\mathbf{a}_t^w \in \mathcal{A}$) based on their local observations \mathbf{o}_t^i as

$$\begin{aligned} \mathbf{g}_t^{m \rightarrow s} &= \pi_m(\mathbf{o}_t^{m \rightarrow s}), & \mathbf{o}_t^{m \rightarrow s} &= \mathbf{h}_t^m \parallel \mathbf{h}_t^{s,0} \parallel \mathbf{h}_t^{s,L_r} & (\text{manager}) \\ \mathbf{g}_t^{s \rightarrow w} &= \pi_s(\mathbf{o}_t^{s \rightarrow w}), & \mathbf{o}_t^{s \rightarrow w} &= \mathbf{g}_t^{m \rightarrow s} \parallel \mathbf{h}_t^{s,L_r} \parallel \mathbf{h}_t^{w,0} \parallel \mathbf{h}_t^{w,L_r} & (\text{sub-manager}) \\ \mathbf{a}_t^w &= \pi_w(\mathbf{o}_t^w), & \mathbf{o}_t^w &= \mathbf{g}_t^{s \rightarrow w} \parallel \mathbf{h}_t^{w,0} \parallel \mathbf{h}_t^{w,L_r} & (\text{worker}) \end{aligned} \quad (2)$$

where π_m, π_s, π_w indicate the policies of the manager, sub-managers, and workers, respectively. Note that each policy effectively acts in a (learned) latent Partially-Observable Markov Decision Process (POMDP) [37], with observations and actions as in Eq. 2 and rewards defined in the

next subsection. As such, value $V^{\pi_i}(\mathbf{o}_t^i)$ and advantage $A^{\pi_i}(\mathbf{o}_t^i, \pi_i(\mathbf{o}_t^i))$ functions associated with observations and actions (goals) can be defined accordingly. Furthermore, upper levels operate at different time scales: we assume that the top-level manager and sub-managers send goals every α_m and α_s environment step, respectively; temporal order in the hierarchy must be preserved [18, 54], i.e., $\alpha_m \geq \alpha_s \geq 1$. Furthermore, to account for goal changes in the most efficient way, it is reasonable to set $\alpha_m = K\alpha_s \doteq K\alpha$, with $K \in \mathbb{Z}^+$. Similarly to other works on HRL [69, 70, 54], K and α are hyperparameters that can be easily tuned according to the problem at hand.

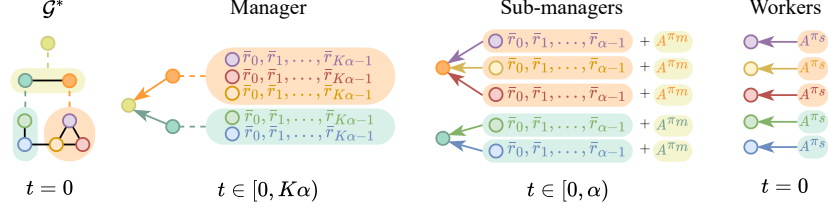


Figure 1: One-step rewards for each level (up to scaling) of the hierarchical graph \mathcal{G}^* at time $t = 0$.

3.2 Building the Hierarchical Graph

The graph \mathcal{G}_t^* encoding the feudal hierarchy is built according to the environment structure. Following the Markovian assumption, it can be defined as a function of the global state \mathbf{S}_t . In a 3-level hierarchy, all the sub-managers are trivially connected to a single top-level manager. The feudal relationships between workers and sub-managers can be derived from a subset of the observation features (e.g., agents' locations) or structural constraints (e.g., communication channels). We found agents' positions to provide a good enough bias to generate the hierarchical and relational structure without requiring additional knowledge on the environment (see Sec. 5). For the upper levels of the hierarchy, we use fully-connected graphs to exchange messages, i.e., we assume that all sub-managers can communicate with each other. Connecting all the workers to a single top-level manager, i.e., using a 2-level hierarchy, is a viable solution for problems with no clear intermediate clustering (see Sec. 5). Hierarchical relationships can also, in principle, be learned end-to-end [45, 21], e.g., by relying on graph pooling operators [28]. While this approach is more general, it comes with overheads in terms of both computational and sample complexity. To keep the scope of the paper contained, we do not explore this direction here.

3.3 Hierarchical Assignment of Rewards

In HiMPo, (sub-)managers operate at different time scales, and the hierarchical graph \mathcal{G}_t^* can be dynamic. Rewards must be assigned locally according to the topology of \mathcal{G}_t^* to avoid sub-optimal solutions. If misspecified, rewards for intermediate levels can lead to undesirable behaviours, such as sub-managers competing against each other. Fig. 1 summarizes the reward-assignment scheme.

Manager Each goal sent by the top-level manager is assigned a local learning signal, based on the rewards obtained by the subset of workers supervised by the sub-manager receiving the goal. In particular, the reward associated with the goal $g_t^{m \rightarrow s}$ sent by the manager m to the s -th sub-manager at time step t reads:

$$r_t^{m \rightarrow s} = \text{MEAN}_{w \in \mathcal{C}_t^s} \left\{ \sum_{k=0}^{K\alpha-1} \bar{r}_{t+k}^w \right\}, \quad (3)$$

where \mathcal{C}_t^s denotes the partition of workers associated with sub-manager s at step t and $\{\bar{r}_k^w\}_{w=1}^N$ is the set of environment rewards at the k -th time step. In other words, $r_t^{m \rightarrow s}$ is defined as the average cumulative reward obtained by agents in \mathcal{C}_t^s over the $K\alpha$ time steps associated with the goal $g_t^{m \rightarrow s}$. Notice that, for the duration of goal $g_t^{m \rightarrow s}$, changes in the subset of workers (and topology \mathcal{G}_t^*) are not observed by the manager m . Therefore, the cumulative reward $r_t^{m \rightarrow s}$ is computed by considering only the workers in \mathcal{C}_t^s , i.e., subordinate workers at the time step when the goal was sent. We compute average (rather than total) cumulative rewards to account for variations in the number of subordinate workers w.r.t. future time steps.

Sub-managers Learning signals for lower levels must be aligned with the objectives set by the upper levels. We define the reward associated with each sub-manager s according to the performance that manager m achieved by assigning the goal $g_t^{m \rightarrow s}$ to s . In particular, the reward signals are defined as the manager’s advantage function, scaled by the number K of sub-manager actions over the $K\alpha$ interval. Maximizing the manager’s advantage function implies taking actions that lead to higher returns, encouraging sub-managers to conform to the set goals. Note that the advantage function $A^{\pi_m}(\mathbf{o}_t^{m \rightarrow s}, \mathbf{g}_t^{m \rightarrow s})$ is conditioned on the goal $g_t^{m \rightarrow s}$ sent by the top-level manager. However, $g_t^{m \rightarrow s}$ is evaluated by aggregating cumulative rewards among workers in \mathcal{C}_t^s (see also Fig. 1). In particular, from the perspective of the sub-manager s , $A^{\pi_m}(\mathbf{o}_t^{m \rightarrow s}, \mathbf{g}_t^{m \rightarrow s})$ is the same regardless of the number of times K in which s acts in the $K\alpha$ steps and the number $|\mathcal{C}_t^s|$ of goals sent at each of the step. Thus, we add a local component to the reward associated with each goal $g_t^{s \rightarrow w}$ to distinguish it from goals assigned to other workers or at different time steps. Specifically, we use the cumulative reward received by worker w over the subsequent α steps. As a result, the reward associated to $g_t^{s \rightarrow w}$ reads:

$$r_t^{s \rightarrow w} = \frac{A^{\pi_m}(\mathbf{o}_t^{m \rightarrow s}, \mathbf{g}_t^{m \rightarrow s})}{K} + \sum_{k=0}^{\alpha-1} \bar{r}_{t+k}^w. \quad (4)$$

Workers We define rewards for workers similarly to sub-managers. In particular, each worker w receives a reward consisting of the advantage function $A^{\pi_s}(\mathbf{o}_t^{s \rightarrow w}, \mathbf{g}_t^{s \rightarrow w})$ of the sub-manager s from which the goal $g_t^{s \rightarrow w}$ has been sent, scaled by the duration α , i.e., as

$$r_t^w = \frac{A^{\pi_s}(\mathbf{o}_t^{s \rightarrow w}, \mathbf{g}_t^{s \rightarrow w})}{\alpha}, \quad (5)$$

where we denote the worker’s reward as r_t^w to distinguish it from the external reward \bar{r}_t^w . For each worker w , the reward in Eq. 5 is distributed among the α actions taken in the interval $[t, t + \alpha)$. Rewards in Eq. 5 are local, i.e., each worker receives a different learning signal. Unlike the sub-managers’ case (Eq. 4), there is no need for adding the external reward here. As a result, workers do not have direct access to external rewards. For dynamic hierarchies, we found that using a modified advantage-like definition of $A^{\pi_s}(\mathbf{o}_t^{s \rightarrow w}, \mathbf{g}_t^{s \rightarrow w})$ yields better results for policy optimization (see App. B.2).

End-to-end Training Following the FRL paradigm, each level-specific policy is trained independently from the others to maximize its level-specific return. This allows for generating goals at different spatial and temporal scales, enabling task decomposition within the hierarchy (see Sec. 5). Although these learning signals depend on the advantage function at upper levels, ensuring they are aligned with each agent’s objective is a significant challenge: we provide theoretical guarantees on this aspect in Sec. 3.4.

Cooperative and Mixed Settings Note that hierarchies with more than 2-levels assume that the environment is cooperative, while a 2-level hierarchy allows for mixed settings. Indeed, in a 3-level hierarchy, the top-level manager aggregates workers’ rewards from the partitions induced by the hierarchical structure (refer to Eq. 3). Conversely, in 2-level hierarchies, where workers are all connected to a single top-level manager, each worker’s return is maximized separately. In this setting, the learning signal of each worker is directly given by the advantage function of the manager, and there is no need to add the local rewards (refer to the second term on the right-hand side of Eq. 4).

3.4 Theoretical Analysis

HiMPo aims to learn a hierarchy of independent policies where each node receives its level-specific reward. To ensure the proposed reward scheme is sound, we must show that each level-specific objective aligns with the global task. Let $T_m = \{0, K\alpha, 2K\alpha, \dots, \infty\}$ denote the manager time scale, i.e., the set of time steps $t \in T_m$ at which each goal $g_t^{m \rightarrow s}$ is sent. Given the joint policy $\pi \doteq (\pi_w, \pi_s, \pi_m)$, the global objective consists of maximizing the return of each worker (agent):

$$\eta(\pi) = \mathbb{E}_{\tau \sim \pi} \left[\sum_{t \in T_m} \gamma^{t/K\alpha} \sum_{k=0}^{K\alpha-1} \bar{r}_{t+k}^w \right], \quad (6)$$

where τ is a trajectory and $\gamma \doteq \gamma_m$ is the discount factor of the top-level manager. The following theoretical results show that maximizing returns obtained from the rewards defined in Sec. 3.3 implies

maximizing the return in Eq. 6. In particular, we assume that each level-specific policy is learned separately while keeping the others fixed. Proofs for all the theoretical results are provided in App. A.

Theorem 3.1 (Manager). *The maximization of the return $\eta_m(\pi_m)$ defined using the rewards in Eq. 3 obtained by manager m sending a goal to sub-manager s implies the maximization of $\eta(\pi)$, i.e.:*

$$\max_{\pi_m} \eta_m(\pi_m) \doteq \max_{\pi_m} \mathbb{E}_{\tau_m \sim \pi} \left[\sum_{t \in T_m} \gamma^{t/K\alpha} r_t^{m \rightarrow s} \right] = \max_{\pi_m} \eta(\pi) \quad (7)$$

It follows that the manager’s objective $\eta_m(\pi_m)$ is aligned with the global task in Eq. 6.

With a slight abuse of notation, let $\hat{\pi}_w$ and $\hat{\pi}_s$ denote the joint policies $(\pi_w, \tilde{\pi}_s, \tilde{\pi}_m)$ and $(\tilde{\pi}_w, \pi_s, \tilde{\pi}_m)$, where $\tilde{\pi}_{i=w,s,m}$ is the old version of $\pi_{i=w,s,m}$. Furthermore, let $T_s = \{0, \alpha, 2\alpha, \dots, \infty\}$ be the sub-managers’ time scale, i.e., the set of time steps t in which goals $g_t^{s \rightarrow w}$ are propagated.

Theorem 3.2 (Sub-managers). *Assuming $\gamma_s \simeq \gamma \simeq 1$ and K not too large, the maximization of the return $\eta_s(\pi_s)$ defined using the rewards in Eq. 4 obtained by the sub-manager s sending a goal to the worker w implies the maximization of $\eta(\hat{\pi}_s)$, i.e.:*

$$\max_{\pi_s} \eta_s(\pi_s) \doteq \max_{\pi_s} \mathbb{E}_{\tau_s \sim \hat{\pi}_s} \left[\sum_{t \in T_s} \gamma_s^{t/\alpha} r_t^{s \rightarrow w} \right] = \max_{\pi_s} \eta(\hat{\pi}_s) \quad (8)$$

It follows that the sub-manager’s objective $\eta_s(\pi_s)$ is aligned with the global objective in Eq. 6.

Theorem 3.3 (Workers). *Assuming $\gamma_w \simeq \gamma_s \simeq 1$ and α not too large, the maximization of the return $\eta_w(\pi_w)$ defined using the rewards in Eq. 5 obtained by the worker w implies the maximization of $\eta(\hat{\pi}_w)$, i.e.:*

$$\max_{\pi_w} \eta_w(\pi_w) \doteq \max_{\pi_w} \mathbb{E}_{\tau_w \sim \hat{\pi}_w} \left[\sum_{t=0}^{\infty} \gamma_w^t r_t^w \right] = \max_{\pi_w} \eta(\hat{\pi}_w) \quad (9)$$

It follows that the worker’s objective $\eta_w(\pi_w)$ is aligned with the global objective in Eq. 6.

Note that the theoretical results remain valid when using a 2-level hierarchy. In such a case, the learning signal of each worker consists of the manager’s advantage function. Finally, as shown in previous work [46], concurrently optimizing policies across different levels – despite assumptions made earlier – can improve sample efficiency without compromising performance.

4 Related Work

Within communication-based methodologies [76], message-passing architectures have been employed to parameterize policies [1, 57] and enhance global awareness during training and execution phases without relying on global information [56]. Other works follow a similar approach in the context of multi-robot systems [38, 14, 13]. GNNs have also been used to learn action-value functions in different settings. Jiang et al. [35] leverage an attention mechanism based on the agent graph to learn individual action-value functions, while other works [55, 41] apply GNNs to learn a global value function in fully cooperative scenarios. Other Comm-MADRL approaches either broadcast messages through static communication structures [24, 66, 17] or directly learn the propagation scheme [65, 34, 47, 32]. However, to the best of our knowledge, none of the previous methods pair message-passing architectures with HRL approaches. Concerning single-agent HRL settings, Li et al. [46] addressed the issue of designing a problem-free reward function for low-level policies by leveraging the advantage function of high-level policies. Subsequently, this reward scheme has been extended to multi-agent systems [70, 49]: here, architectures are trained under the CTDE paradigm using QMIX [61]. Different from our approach, the optimization signal of low-level policies is based on the team reward, i.e., actions are not explicitly evaluated according to their local contributions [67]. Furthermore, such architectures are restricted to 2-level hierarchies. The FRL [18] framework has been first extended to the deep RL setting [69] and, subsequently, applied to multi-agent systems [2]. Within this context, it showed performance improvements across different domains, ranging from traffic signal control [51, 52] to warehouse logistics [42]. Nevertheless, none of these multi-agent methods adopt a problem-free reward scheme nor consider message-passing policies. More recently, the FRL paradigm has been exploited to learn hierarchical graph-based modular policies in single-agent RL [54].

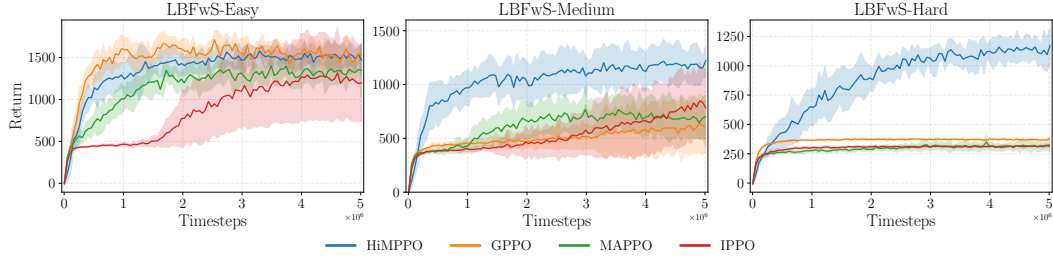


Figure 2: Results of the models in the *LBFwS* environment. For each configuration, we report the average return and sample standard deviation of 8 runs.

5 Experiments

In the following, we evaluate HiMPo on relevant MARL benchmarks by considering problems of increasing complexity and with a varying number of interacting agents. We use Proximal Policy Optimization (PPO) [64] as the underlying policy optimization algorithm and refer to this specific implementation of HiMPo as Hierarchical Message-Passing Proximal Policy Optimization (HiMPPO). As anticipated, we optimize all the levels concurrently, as it has been shown to improve sample efficiency [46].

Baselines and benchmarks We compare our methodology against state-of-the-art MARL methods. In particular, we consider 1) *IPPO* and 2) *MAPPO*, i.e., two MARL variants of PPO proposed by Yu et al. [73]. Concerning graph-based methods, we include in the comparison 3) the homogeneous version of Graph Neural Network PPO (*GPPO*) [13]. We consider 2 main environments under several different settings: 1) a novel version of the *Level-Based Foraging* (LBF) [3, 59] environment, to show how the models perform in tasks with challenging temporal credit assignment; 2) the *VMAS Sampling* environment from the VMAS simulator [12], to assess the exploration capabilities of the baselines. We report results on 3) *SMACv2* [22], i.e., the improved version of the *StarCraft Multi-Agent Challenge* (SMAC) [62], as an additional benchmark in App. B.1. For the models using graph-based representations, we use a dynamic graph based on the location of the agents at each time step. In particular, two agents are connected if they are within their observation ranges. For *LBFwS* and *VMAS Sampling* we use the standard observation ranges, while for *SMACv2*, we set it to 3. Moreover, in *SMACv2*, dead allies are filtered out from the graph. All the models were tuned on *LBFwS*, and the same sets of hyperparameters were used for the other environments. Additional details regarding environments and models are reported in App. C.1 and C.3.

5.1 Level-Based Foraging with Survival (LBFwS)

In *LBFwS* (App. C.1.1), agents navigate in a grid world where resources of different value are randomly spawned. Consuming resources increases the individual reward, while delivering them to a fixed location extends the episode duration. Higher-value items need the cooperation of multiple agents to be obtained. An agent can learn either to simply maximize its reward (individual strategy) or to sacrifice the local reward to extend the episode and possibly achieve a higher final return (cooperative strategy). In particular, we consider a system of 10 agents and three difficulty levels (*Easy*, *Medium*, and *Hard*) that correspond to progressively larger grid sizes and varying amount of resources (see App. C.1.1). For HiMPPO, we use a static 2-level hierarchical graph \mathcal{G}^* with all the agents (workers) connected to a single top-level manager that sends goals every 5 steps.

Results Results are reported in Fig. 2. All the models can achieve high returns in the *LBFwS-Easy* scenario, with the graph-based methods (HiMPPO and GPPO) being more sample efficient. However, as the difficulty increases, the baselines struggle to achieve good results. Notably, in *LBFwS-Hard*, all the methods except HiMPPO learn to maximize the individual reward. Conversely, HiMPPO consistently achieves good results at every difficulty level. These results show the effectiveness of HiMPo in coordination and planning capabilities.

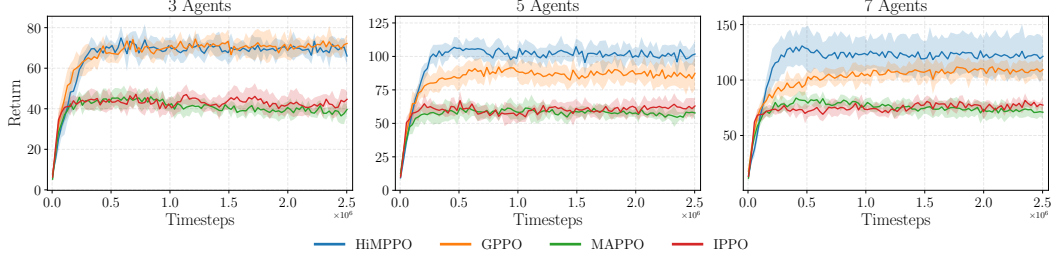


Figure 3: Results of the models for different numbers of agents in the *VMAS Sampling* environment. For each configuration, we report the average return and sample standard deviation of 6 runs.

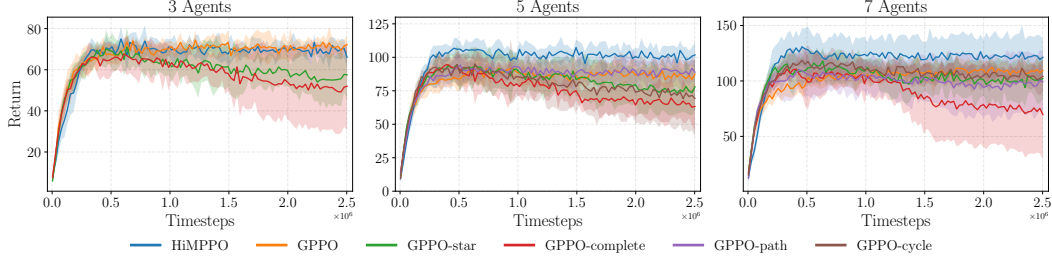


Figure 4: Analysis of different communication topologies in the *VMAS Sampling* environment. For each configuration, we report the average return and sample standard deviation of 6 runs.

5.2 VMAS Sampling

In this continuous control problem, n robots are randomly spawned in a grid. The grid has an underlying Gaussian density function with 3 modes. Visiting a cell for the first time leads to a reward proportional to the density function. For HiMPPO, the dynamic 3-level hierarchical graph \mathcal{G}_t^* is generated according to the positions of the robots at each step t (refer to App. C.1.2); the goal propagation scales are set to $K = 2$ and $\alpha = 5$. In the advantage-like definition of the workers’ rewards (App. A.4), future signals are defined using a one-step truncation on the sequence of sub-managers’ value functions (refer to Eq. 19), as empirical evidence showed better performance. App. B.2 further analyzes the impact of considering different credit assignment schemes in this environment.

Results The results in Fig. 3 show that HiMPPO performs better than the baselines. As expected, graph-based approaches (HiMPPO and GPPO) outperform the baselines that do not implement communication, as they can effectively share valuable information during both training and execution. Nevertheless, HiMPPO achieves higher overall returns when dealing with larger systems, showing the benefits of relying on hierarchical policies to coordinate the exploration.

5.3 Analysis and Ablations

In this section, we perform a set of ablation studies and additional analyses to assess the impact and effectiveness of the proposed designs.

Communication Mechanism To analyze the impact of the proposed coordination and communication mechanism, we compare the results achieved by HiMPPO against some variants of GPPO that use a fixed graph topology. In particular, we consider a star graph (GPPO-star), with a central node connected to all the others, a fully connected graph (GPPO-complete), a path graph (GPPO-path), and a cycle graph (GPPO-cycle). Note that for 3-node graphs, path and cycle graphs are equivalent to star and fully-connected graphs, respectively. Results reported in Fig. 4 show that HiMPPO performs better than all the variants for larger systems. In particular, using fixed topologies in GPPO often hinders performance, and using a fully-connected graph can make optimization unstable. Note that globally HiMPPO and fixed-topology variants of GPPO have access to the same information, showing that *how* such information is processed is a crucial aspect in learning effective policies.

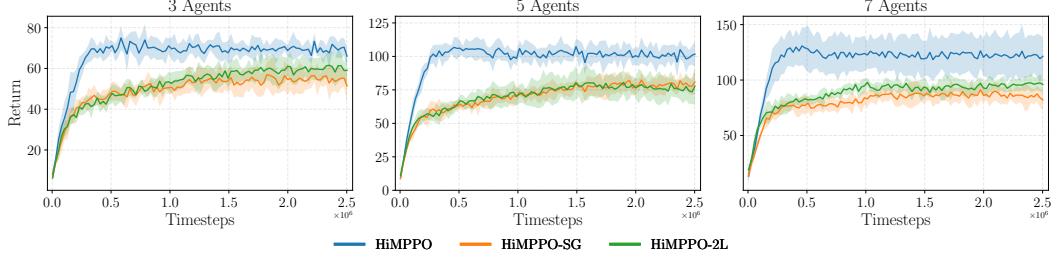


Figure 5: Ablation study for different hierarchical graphs in the *VMAS Sampling* environment. For each configuration, we report the average return and sample standard deviation of 6 runs.

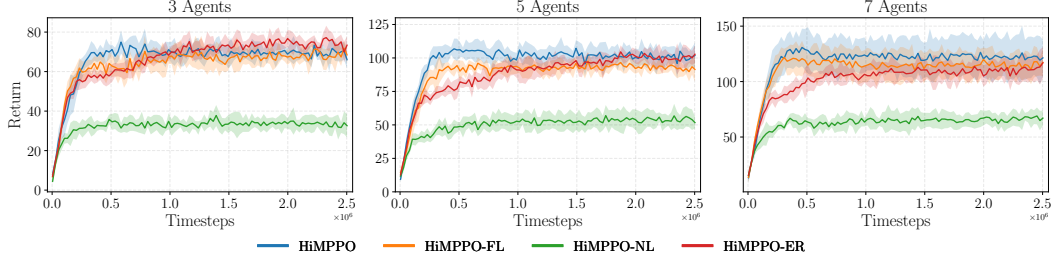


Figure 6: Ablation study for different reward-assignment schemes in the *VMAS Sampling* environment. For each configuration, we report the average return and sample standard deviation of 6 runs.

Hierarchical Structure We consider different hierarchical graph structures to investigate their impact on decision making. In particular, we compare the state-dependent dynamic hierarchical graph \mathcal{G}_t^* (refer to Sec. 3.2 and App. C.1.2) against: 1) Static Graph HiMPPO (HiMPPO-SG), a 3-level HiMPPO where the topology is fixed at the beginning of each episode, i.e., the $\mathcal{G}_t^* = \mathcal{G}_0^* \forall t$; 2) 2-Levels HiMPPO (HiMPPO-2L), where we use a static 2-level hierarchy with all the workers connected to a single top-level manager. The results reported in Fig. 5 show that the dynamic 3-level hierarchy is crucial for achieving better performance in *VMAS Sampling*. Indeed, HiMPPO-SG and HiMPPO-2L are limited compared to HiMPPO since the hierarchical structure does not reflect the current state of the environment. Notably, HiMPPO-SG and HiMPPO-2L perform akin. This suggests that, in this environment, processing information with a deeper (but static) hierarchy does not provide additional advantages when compared to a 2-level structure.

Reward Assignment To validate the hierarchical reward-assignment scheme proposed in Sec. 3.3, we compare HiMPPO against different variants. In particular, we consider: 1) Full Local HiMPPO (HiMPPO-FL), where we expose workers to the external reward by adding it to Eq. 5; 2) No Local HiMPPO (HiMPPO-NL), where sub-managers do not receive the cumulative reward associated with each worker (refer to Eq. 4); 3) Environment Reward HiMPPO (HiMPPO-ER), where rewards of sub-managers and workers are based only on the external signal (the top-level manager still receives the reward in Eq. 3). The results reported in Fig. 6 support adopting the proposed hierarchical reward-assignment scheme. In particular, the comparison with HiMPPO-FL and HiMPPO-ER shows that preventing workers from accessing the external signal and employing advantage-like rewards does not deteriorate performance and improves sample efficiency. Conversely, preventing sub-managers from using local rewards leads to poor results, as shown by HiMPPO-NL.

6 Conclusions

We introduced *Hierarchical Message-Passing Policy* (HiMPo), a novel methodology for learning hierarchies of message-passing policies in multi-agent systems. The decision-making structure of HiMPo relies on a (dynamic) multi-level hierarchical graph, which allows for improved coordination and planning by leveraging communication and feudal relationships. Notably, the proposed reward-assignment method is theoretically sound, and empirical results on different benchmarks validate its effectiveness. HiMPo opens up new designs for graph-based hierarchical MARL policies.

Limitations and Future Works The current implementation of HiMPo employs a goal-assignment scheme where upper levels operate at fixed time scales. Future directions might build upon HiMPo to explore more sophisticated feudal mechanisms and address these limitations. In this sense, it would be interesting to develop an asynchronous and adaptive goal-assignment mechanism, where each (sub-)manager can propagate commands at different time scales according to the current state. Finally, we believe that extending the methodology to heterogeneous settings with more advanced hierarchical dependencies could broaden its applicability.

7 Acknowledgements

This work was supported by the Swiss National Science Foundation grants No. 204061 (*HORD GNN: Higher-Order Relations and Dynamics in Graph Neural Networks*) and No. 225351 (*Relational Deep Learning for Reliable Time Series Forecasting at Scale*).

References

- [1] Agarwal, A., Kumar, S., Sycara, K., and Lewis, M. Learning transferable cooperative behavior in multi-agent teams. In *Proceedings of the 19th International Conference on Autonomous Agents and MultiAgent Systems*, AAMAS '20, pp. 1741–1743, Richland, SC, 2020. International Foundation for Autonomous Agents and Multiagent Systems. ISBN 9781450375184.
- [2] Ahilan, S. and Dayan, P. Feudal multi-agent hierarchies for cooperative reinforcement learning. *arXiv preprint arXiv:1901.08492*, 2019.
- [3] Albrecht, S. V. and Ramamoorthy, S. A game-theoretic model and best-response learning method for ad hoc coordination in multiagent systems. In *Proceedings of the 2013 International Conference on Autonomous Agents and Multi-Agent Systems*, AAMAS '13, pp. 1155–1156, Richland, SC, 2013. International Foundation for Autonomous Agents and Multiagent Systems. ISBN 9781450319935.
- [4] Albrecht, S. V., Christianos, F., and Schäfer, L. *Multi-Agent Reinforcement Learning: Foundations and Modern Approaches*. MIT Press, 2024. URL <https://www.marl-book.com>.
- [5] Aloor, J. J., Nayak, S., Dolan, S., Qin, V. L., and Balakrishnan, H. Towards cooperation and fairness in multi-agent reinforcement learning. In *Coordination and Cooperation for Multi-Agent Reinforcement Learning Methods Workshop*, 2024. URL <https://openreview.net/forum?id=ze0wbBNQrA>.
- [6] Amato, C. An introduction to centralized training for decentralized execution in cooperative multi-agent reinforcement learning. *arXiv preprint arXiv:2409.03052*, 2024.
- [7] Azmani, H., Rosseau, A., Nowe, A., and Rădulescu, R. Cooperative foraging behaviour through multi-agent reinforcement learning with graph-based communication. In *Sixteenth European Workshop on Reinforcement Learning*, 2023.
- [8] Bacciu, D., Errica, F., Micheli, A., and Podda, M. A gentle introduction to deep learning for graphs. *Neural Networks*, 129:203–221, 2020.
- [9] Barhate, N. Minimal pytorch implementation of proximal policy optimization. <https://github.com/nikhilbarhate99/PP0-PyTorch>, 2021.
- [10] Barto, A. G. and Mahadevan, S. Recent advances in hierarchical reinforcement learning. *Discrete event dynamic systems*, 13(1):41–77, 2003.
- [11] Battaglia, P. W., Hamrick, J. B., Bapst, V., Sanchez-Gonzalez, A., Zambaldi, V., Malinowski, M., Tacchetti, A., Raposo, D., Santoro, A., Faulkner, R., et al. Relational inductive biases, deep learning, and graph networks. *arXiv preprint arXiv:1806.01261*, 2018.
- [12] Bettini, M., Kortvelesy, R., Blumenkamp, J., and Prorok, A. Vmas: A vectorized multi-agent simulator for collective robot learning. *The 16th International Symposium on Distributed Autonomous Robotic Systems*, 2022.
- [13] Bettini, M., Shankar, A., and Prorok, A. Heterogeneous multi-robot reinforcement learning. In *Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent Systems*, pp. 1485–1494, 2023.
- [14] Blumenkamp, J., Morad, S., Gielis, J., Li, Q., and Prorok, A. A framework for real-world multi-robot systems running decentralized gnn-based policies. In *2022 International Conference on Robotics and Automation (ICRA)*, pp. 8772–8778. IEEE, 2022.
- [15] Bou, A., Bettini, M., Dittert, S., Kumar, V., Sodhani, S., Yang, X., Fabritiis, G. D., and Moens, V. Torchrl: A data-driven decision-making library for pytorch, 2023.
- [16] Bronstein, M. M., Bruna, J., Cohen, T., and Veličković, P. Geometric deep learning: Grids, groups, graphs, geodesics, and gauges. *arXiv preprint arXiv:2104.13478*, 2021.
- [17] Das, A., Gervet, T., Romoff, J., Batra, D., Parikh, D., Rabbat, M., and Pineau, J. Tarmac: Targeted multi-agent communication. In *International Conference on machine learning*, pp. 1538–1546. PMLR, 2019.

- [18] Dayan, P. and Hinton, G. E. Feudal reinforcement learning. *Advances in neural information processing systems*, 5, 1992.
- [19] De Witt, C. S., Gupta, T., Makoviichuk, D., Makoviyhuk, V., Torr, P. H., Sun, M., and Whiteson, S. Is independent learning all you need in the starcraft multi-agent challenge? *arXiv preprint arXiv:2011.09533*, 2020.
- [20] Duan, W., Lu, J., and Xuan, J. Inferring latent temporal sparse coordination graph for multi-agent reinforcement learning. *CoRR*, abs/2403.19253, 2024.
- [21] Duan, W., Lu, J., and Xuan, J. Group-aware coordination graph for multi-agent reinforcement learning. In Larson, K. (ed.), *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence, IJCAI-24*, pp. 3926–3934. International Joint Conferences on Artificial Intelligence Organization, 8 2024. doi: 10.24963/ijcai.2024/434. URL <https://doi.org/10.24963/ijcai.2024/434>. Main Track.
- [22] Ellis, B., Cook, J., Moalla, S., Samvelyan, M., Sun, M., Mahajan, A., Foerster, J., and Whiteson, S. Smacv2: An improved benchmark for cooperative multi-agent reinforcement learning. *Advances in Neural Information Processing Systems*, 36:37567–37593, 2023.
- [23] Fey, M. and Lenssen, J. E. Fast graph representation learning with pytorch geometric. *arXiv preprint arXiv:1903.02428*, 2019.
- [24] Foerster, J., Assael, I. A., De Freitas, N., and Whiteson, S. Learning to communicate with deep multi-agent reinforcement learning. *Advances in neural information processing systems*, 29, 2016.
- [25] Ghavamzadeh, M., Mahadevan, S., and Makar, R. Hierarchical multi-agent reinforcement learning. *Autonomous Agents and Multi-Agent Systems*, 13:197–229, 2006.
- [26] Gilmer, J., Schoenholz, S. S., Riley, P. F., Vinyals, O., and Dahl, G. E. Neural message passing for quantum chemistry. In *International conference on machine learning*, pp. 1263–1272. PMLR, 2017.
- [27] Grammel, N., Son, S., Black, B., and Agrawal, A. Revisiting parameter sharing in multi-agent deep reinforcement learning. *arXiv preprint arXiv:2005.13625*, 2020.
- [28] Grattarola, D., Zambon, D., Bianchi, F. M., and Alippi, C. Understanding pooling in graph neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- [29] Gupta, J. K., Egorov, M., and Kochenderfer, M. J. Cooperative multi-agent control using deep reinforcement learning. In *AAMAS Workshops*, 2017. URL <https://api.semanticscholar.org/CorpusID:9421360>.
- [30] Harris, C. R., Millman, K. J., Van Der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau, D., Wieser, E., Taylor, J., Berg, S., Smith, N. J., et al. Array programming with numpy. *Nature*, 585 (7825):357–362, 2020.
- [31] Hernandez-Leal, P., Kaisers, M., Baarslag, T., and De Cote, E. M. A survey of learning in multiagent environments: Dealing with non-stationarity. *arXiv preprint arXiv:1707.09183*, 2017.
- [32] Hu, S., Shen, L., Zhang, Y., and Tao, D. Learning multi-agent communication from graph modeling perspective. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=Qox9r00kN0>.
- [33] Huh, D. and Mohapatra, P. Multi-agent reinforcement learning: A comprehensive survey. *arXiv preprint arXiv:2312.10256*, 2023.
- [34] Jiang, J. and Lu, Z. Learning attentional communication for multi-agent cooperation. *Advances in neural information processing systems*, 31, 2018.
- [35] Jiang, J., Dun, C., Huang, T., and Lu, Z. Graph convolutional reinforcement learning. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=HkxdQkSYDB>.

- [36] Jin, W., Du, H., Zhao, B., Tian, X., Shi, B., and Yang, G. A comprehensive survey on multi-agent cooperative decision-making: Scenarios, approaches, challenges and perspectives, 2025. URL <https://arxiv.org/abs/2503.13415>.
- [37] Kaelbling, L. P., Littman, M. L., and Cassandra, A. R. Planning and acting in partially observable stochastic domains. *Artificial intelligence*, 101(1-2):99–134, 1998.
- [38] Khan, A., Tolstaya, E., Ribeiro, A., and Kumar, V. Graph policy gradients for large scale robot control. In *Conference on robot learning*, pp. 823–834. PMLR, 2020.
- [39] Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [40] Kipf, T. N. and Welling, M. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations*, 2017. URL <https://openreview.net/forum?id=SJU4ayYgl>.
- [41] Kortvelesy, R. and Prorok, A. Qgnn: Value function factorisation with graph neural networks. *arXiv preprint arXiv:2205.13005*, 2022.
- [42] Krnjaic, A., Steleac, R. D., Thomas, J. D., Papoudakis, G., Schäfer, L., To, A. W. K., Lao, K.-H., Cubuktepe, M., Haley, M., Börsting, P., et al. Scalable multi-agent reinforcement learning for warehouse logistics with robotic and human co-workers. *arXiv preprint arXiv:2212.11498*, 2022.
- [43] Kulkarni, T. D., Narasimhan, K., Saeedi, A., and Tenenbaum, J. Hierarchical deep reinforcement learning: Integrating temporal abstraction and intrinsic motivation. *Advances in neural information processing systems*, 29, 2016.
- [44] Kurin, V., Igl, M., Rocktäschel, T., Boehmer, W., and Whiteson, S. My body is a cage: the role of morphology in graph-based incompatible control. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=N3zUDGN510>.
- [45] Li, J., Hua, C., Park, J., Ma, H., Dax, V. M., and Kochenderfer, M. J. Evolvehypergraph: Group-aware dynamic relational reasoning for trajectory prediction. *CoRR*, abs/2208.05470, 2022. URL <https://doi.org/10.48550/arXiv.2208.05470>.
- [46] Li, S., Wang, R., Tang, M., and Zhang, C. Hierarchical reinforcement learning with advantage-based auxiliary rewards. *Advances in Neural Information Processing Systems*, 32, 2019.
- [47] Li, S., Gupta, J. K., Morales, P., Allen, R., and Kochenderfer, M. J. Deep implicit coordination graphs for multi-agent reinforcement learning. In *Proceedings of the 20th International Conference on Autonomous Agents and MultiAgent Systems*, pp. 764–772, 2021.
- [48] Littman, M. L. Markov games as a framework for multi-agent reinforcement learning. In *Machine learning proceedings 1994*, pp. 157–163. Elsevier, 1994.
- [49] Liu, Z., Xu, Z., and Fan, G. Hierarchical multi-agent reinforcement learning with intrinsic reward rectification. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1–5. IEEE, 2023.
- [50] Lyu, X., Baisero, A., Xiao, Y., Daley, B., and Amato, C. On centralized critics in multi-agent reinforcement learning. *Journal of Artificial Intelligence Research*, 77:295–354, 2023.
- [51] Ma, J. and Wu, F. Feudal multi-agent deep reinforcement learning for traffic signal control. In *Proceedings of the 19th international conference on autonomous agents and multiagent systems (AAMAS)*, pp. 816–824, 2020.
- [52] Ma, J. and Wu, F. Feudal multi-agent reinforcement learning with adaptive network partition for traffic signal control. *arXiv preprint arXiv:2205.13836*, 2022.
- [53] Makar, R., Mahadevan, S., and Ghavamzadeh, M. Hierarchical multi-agent reinforcement learning. In *Proceedings of the fifth international conference on Autonomous agents*, pp. 246–253, 2001.

- [54] Marzi, T., Khehra, A. S., Cini, A., and Alippi, C. Feudal graph reinforcement learning. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL <https://openreview.net/forum?id=wFcyJTik90>.
- [55] Naderializadeh, N., Hung, F. H., Soleyman, S., and Khosla, D. Graph convolutional value decomposition in multi-agent reinforcement learning. *arXiv preprint arXiv:2010.04740*, 2020.
- [56] Nayak, S., Choi, K., Ding, W., Dolan, S., Gopalakrishnan, K., and Balakrishnan, H. Scalable multi-agent reinforcement learning through intelligent information aggregation. In *International Conference on Machine Learning*, pp. 25817–25833. PMLR, 2023.
- [57] Niu, Y., Paleja, R. R., and Gombolay, M. C. Multi-agent graph-attention communication and teaming. In *AAMAS*, volume 21, pp. 20th, 2021.
- [58] Papoudakis, G., Christianos, F., Rahman, A., and Albrecht, S. V. Dealing with non-stationarity in multi-agent deep reinforcement learning. *arXiv preprint arXiv:1906.04737*, 2019.
- [59] Papoudakis, G., Christianos, F., Schäfer, L., and Albrecht, S. V. Benchmarking multi-agent deep reinforcement learning algorithms in cooperative tasks. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*, 2021. URL <https://openreview.net/forum?id=cIrPX-Sn5n>.
- [60] Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Köpf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., and Chintala, S. *PyTorch: an imperative style, high-performance deep learning library*. Curran Associates Inc., Red Hook, NY, USA, 2019.
- [61] Rashid, T., Samvelyan, M., De Witt, C. S., Farquhar, G., Foerster, J., and Whiteson, S. Monotonic value function factorisation for deep multi-agent reinforcement learning. *Journal of Machine Learning Research*, 21(178):1–51, 2020.
- [62] Samvelyan, M., Rashid, T., De Witt, C. S., Farquhar, G., Nardelli, N., Rudner, T. G., Hung, C.-M., Torr, P. H., Foerster, J., and Whiteson, S. The starcraft multi-agent challenge. *arXiv preprint arXiv:1902.04043*, 2019.
- [63] Schulman, J., Moritz, P., Levine, S., Jordan, M., and Abbeel, P. High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015.
- [64] Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms. *CoRR*, abs/1707.06347, 2017. URL <http://arxiv.org/abs/1707.06347>.
- [65] Singh, A., Jain, T., and Sukhbaatar, S. Learning when to communicate at scale in multiagent cooperative and competitive tasks. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=rye7knCqK7>.
- [66] Sukhbaatar, S., Fergus, R., et al. Learning multiagent communication with backpropagation. *Advances in neural information processing systems*, 29, 2016.
- [67] Sunehag, P., Lever, G., Gruslys, A., Czarnecki, W. M., Zambaldi, V., Jaderberg, M., Lanctot, M., Sonnerat, N., Leibo, J. Z., Tuyls, K., and Graepel, T. Value-decomposition networks for cooperative multi-agent learning based on team reward. In *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems, AAMAS ’18*, pp. 2085–2087, Richland, SC, 2018. International Foundation for Autonomous Agents and Multiagent Systems.
- [68] Sutton, R. S., Precup, D., and Singh, S. Between mdps and semi-mdps: A framework for temporal abstraction in reinforcement learning. *Artificial intelligence*, 112(1-2):181–211, 1999.
- [69] Vezhnevets, A. S., Osindero, S., Schaul, T., Heess, N., Jaderberg, M., Silver, D., and Kavukcuoglu, K. Feudal networks for hierarchical reinforcement learning. In *International conference on machine learning*, pp. 3540–3549. PMLR, 2017.
- [70] Xu, Z., Bai, Y., Zhang, B., Li, D., and Fan, G. Haven: hierarchical cooperative multi-agent reinforcement learning with dual coordination mechanism. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 11735–11743, 2023.

- [71] Yadan, O. Hydra - a framework for elegantly configuring complex applications. Github, 2019. URL <https://github.com/facebookresearch/hydra>.
- [72] Yang, J., Borovikov, I., and Zha, H. Hierarchical cooperative multi-agent reinforcement learning with skill discovery. In *Proceedings of the 19th International Conference on Autonomous Agents and MultiAgent Systems*, pp. 1566–1574, 2020.
- [73] Yu, C., Velu, A., Vinitsky, E., Gao, J., Wang, Y., Bayen, A., and Wu, Y. The surprising effectiveness of ppo in cooperative multi-agent games. *Advances in Neural Information Processing Systems*, 35:24611–24624, 2022.
- [74] Yuan, L., Zhang, Z., Li, L., Guan, C., and Yu, Y. A survey of progress on cooperative multi-agent reinforcement learning in open environment. *arXiv preprint arXiv:2312.01058*, 2023.
- [75] Zhou, Y., Liu, S., Qing, Y., Chen, K., Zheng, T., Huang, Y., Song, J., and Song, M. Is centralized training with decentralized execution framework centralized enough for marl? *arXiv preprint arXiv:2305.17352*, 2023.
- [76] Zhu, C., Dastani, M., and Wang, S. A survey of multi-agent deep reinforcement learning with communication. *Autonomous Agents and Multi-Agent Systems*, 38(1):4, 2024.

Appendix

A Additional Proofs

The topology of the hierarchical graph \mathcal{G}_t^* can change in time. In the execution phase, received states are aggregated according to the current topology, making nodes aware of the hierarchical decision-making structure. Nevertheless, long-term returns do not provide information concerning possible changes in \mathcal{G}_t^* . As such, during the training phase, nodes cannot perceive changes in the topology using the sole learning signal. For example, consider the top-level manager: at each step $t \in T_m$, goals $\mathbf{g}_t^{m \rightarrow s}$ are sent according to \mathcal{G}_t^* , which accounts for differences in the partitions of the underneath nodes. However, in the optimization phase, such differences do not emerge when considering the long-term return. To address this problem, during training, we assume a static hierarchy given by the topology at the first step. In other words, each trajectory is evaluated using the hierarchical graph \mathcal{G}_0^* at the first step $t = 0$.

A.1 Proof of Theorem 3.1

Proof. Following the Markovian assumption, the probability of having a subset \mathcal{C}_t^s of workers subordinate to s at time t depends only on the global state \mathbf{S}_t . The top-level manager takes actions according to the current topology \mathcal{G}_t^* . Therefore, long-term returns achieved by sending a goal $\mathbf{g}_t^{m \rightarrow s}$ must be evaluated by considering the rewards associated with the subset of workers \mathcal{C}_s^t , i.e., the objective of the manager is defined by aggregating rewards w.r.t. the initial partition of the workers. This is a reasonable assumption since the scalar learning signal (return) received by the manager does not provide any information concerning possible changes in the topology in the subsequent steps: the manager learns to provide optimal actions according to the current topology \mathcal{G}_t^* , and returns must be coherently defined under this setting. The manager's objective reads:

$$\begin{aligned} \eta_m(\pi_m) &= \mathbb{E}_{\mathbf{S}_0} \mathbb{E}_{\mathbf{o}_0^{m \rightarrow s}} \left[V^{\pi_m}(\mathbf{o}_0^{m \rightarrow s}) \right] = \mathbb{E}_{\tau_m \sim \pi} \left[\sum_{t \in T_m} \gamma^{t/K\alpha} r_t^{m \rightarrow s} \right] = \\ &= \mathbb{E}_{\mathbf{S}_0} \mathbb{E}_{\mathcal{C}_0^s} \mathbb{E}_{\tau_m \sim \pi} \left\{ \sum_{t \in T_m} \gamma^{t/K\alpha} \text{MEAN}_{w \in \mathcal{C}_0^s} \left[\sum_{k=0}^{K\alpha-1} \bar{r}_{t+k}^w \right] \right\} = \\ &= \mathbb{E}_{\mathbf{S}_0} \mathbb{E}_{\mathcal{C}_0^s} \left\{ \text{MEAN}_{w \in \mathcal{C}_0^s} \left[\mathbb{E}_{\tau_m \sim \pi} \left[\sum_{t \in T_m} \gamma^{t/K\alpha} \sum_{k=0}^{K\alpha-1} \bar{r}_{t+k}^w \right] \right] \right\} = \mathbb{E}_{\mathbf{S}_0} \mathbb{E}_{\mathcal{C}_0^s} \left[\text{MEAN}_{w \in \mathcal{C}_0^s} \left\{ \eta(\pi) \right\} \right] = k_m \eta(\pi) \quad (10) \end{aligned}$$

In the last equality, we used the homogeneity of the agents, and $k_m > 0$ is a structural coefficient indicating the average number of subordinate nodes for each sub-manager. Therefore, the manager's objective is aligned with the optimization goal of the joint policy (refer to Eq. 6). This proves Theorem 3.1. \square

A.2 Intermediate Results

In order to provide guarantees of alignment for the lower levels, we now report two intermediate results. To show that low-level returns are aligned with the global objective, let $\tilde{\pi}$ and $\tilde{\pi}_{i=w,s,m}$ denote the old versions of π and π_i , respectively. Consider the following auxiliary term:

$$\eta_{adv}(\pi) \doteq \mathbb{E}_{\tau_m \sim \pi} \left[\sum_{t \in T_m} \gamma^{t/K\alpha} A^{\tilde{\pi}_m}(\mathbf{o}_t^{m \rightarrow s}, \mathbf{g}_t^{m \rightarrow s}) \right] \quad (11)$$

Lemma A.1. *The objective $\eta(\pi)$ can be written in terms of the advantage over the old version $\tilde{\pi}$ of the joint policy as:*

$$\eta(\pi) = \eta(\tilde{\pi}) + \frac{1}{k_m} \eta_{adv}(\pi) \quad (12)$$

Proof. We can write the expected return $\eta(\pi)$ of the policy π in terms of its expected advantage over the old version $\tilde{\pi}$:

$$\begin{aligned}
\mathbb{E}_{\tau_m \sim \pi} \left[\sum_{t \in T_m} \gamma^{t/K\alpha} A^{\tilde{\pi}_m}(\mathbf{o}_t^{m \rightarrow s}, \mathbf{g}_t^{m \rightarrow s}) \right] &= \\
&= \mathbb{E}_{\tau_m \sim \pi} \left[\sum_{t \in T_m} \gamma^{t/K\alpha} [r_t^{m \rightarrow s} + \gamma V^{\tilde{\pi}_m}(\mathbf{o}_{t+K\alpha}^{m \rightarrow s}) - V^{\tilde{\pi}_m}(\mathbf{o}_t^{m \rightarrow s})] \right] = \\
&= \mathbb{E}_{\tau_m \sim \pi} \left[\sum_{t \in T_m} \gamma^{t/K\alpha} r_t^{m \rightarrow s} - V^{\tilde{\pi}_m}(\mathbf{o}_0^{m \rightarrow s}) \right] = \\
&= \mathbb{E}_{\tau_m \sim \pi} \sum_{t \in T_m} \gamma^{t/K\alpha} r_t^{m \rightarrow s} - \mathbb{E}_{\mathbf{s}_0} \mathbb{E}_{\mathbf{o}_0^{m \rightarrow s}} V^{\tilde{\pi}_m}(\mathbf{o}_0^{m \rightarrow s}) = [\eta_m(\pi_m) - \eta_m(\tilde{\pi}_m)] = k_m [\eta(\pi) - \eta(\tilde{\pi})]
\end{aligned} \tag{13}$$

where in the last equality we used [Theorem 3.1](#). Rearranging the terms, we obtain the equality in [Eq. 12](#). \square

Corollary A.2. Given [Lemma A.1](#), maximizing the global objective $\eta(\pi)$ can be written as:

$$\max_{\pi} \eta(\pi) = \max_{\pi} \eta_{adv}(\pi) \tag{14}$$

Proof. Since $\eta(\tilde{\pi})$ is independent of the joint policy π , it is straightforward to see that maximizing $\eta(\pi)$ is independent of $\eta(\tilde{\pi})$ and the optimization of π can be expressed as [Eq. 14](#). \square

A.3 Proof of [Theorem 3.2](#)

Proof. As we did for the top-level manager, the return of the sub-manager s sending goal $\mathbf{g}_t^{s \rightarrow w}$ to the worker w is computed by considering the sequence of rewards obtained by w . We remark that in the $K\alpha$ steps where the top-level manager is not sending goals, the advantage function of the manager in [Eq. 4](#) is fixed for each of the K goals sent by the sub-manager. Therefore, the objective for the sub-manager s reads:

$$\begin{aligned}
\eta_s(\pi_s) &= \mathbb{E}_{\mathbf{s}_0} \mathbb{E}_{\mathbf{o}_0^{s \rightarrow w}} [V^{\hat{\pi}_s}(\mathbf{o}_0^{s \rightarrow w})] = \mathbb{E}_{\tau_s \sim \hat{\pi}_s} \left\{ \sum_{t \in T_s} \gamma_s^{t/\alpha} r_t^{s \rightarrow w} \right\} = \\
&= \mathbb{E}_{\tau \sim \hat{\pi}_s} \left\{ \sum_{t \in T_s} \gamma_s^{t/\alpha} \left[\frac{1}{K} A^{\tilde{\pi}_m}(\mathbf{o}_t^{m \rightarrow s}, \mathbf{g}_t^{m \rightarrow s}) + \sum_{i=0}^{\alpha-1} \bar{r}_{t+i}^w \right] \right\} = \\
&= \mathbb{E}_{\tau_m \sim \hat{\pi}_s} \left\{ \sum_{t \in T_m} \mathbb{E}_{\tau_s(t) \sim \hat{\pi}_s} \left[\sum_{k=0}^{K-1} \mathbb{E}_{\tau_w(k) \sim \hat{\pi}_s} \gamma_s^{t/\alpha} \gamma_s^k \left(\frac{1}{K} A^{\tilde{\pi}_m}(\mathbf{o}_t^{m \rightarrow s}, \mathbf{g}_t^{m \rightarrow s}) + \sum_{i=0}^{\alpha-1} \bar{r}_{t+k\alpha+i}^w \right) \right] \right\} \simeq \\
&\simeq \mathbb{E}_{\tau_m \sim \hat{\pi}_s} \left\{ \sum_{t \in T_m} \frac{\gamma^{t/K\alpha}}{K} \frac{1-\gamma^K}{1-\gamma} A^{\tilde{\pi}_m}(\mathbf{o}_t^{m \rightarrow s}, \mathbf{g}_t^{m \rightarrow s}) \right\} + \mathbb{E}_{\tau_m \sim \hat{\pi}_s} \left[\sum_{t \in T_m} \left(\gamma^{t/K\alpha} \sum_{i=0}^{K\alpha-1} \bar{r}_{t+i}^w \right) \right] = \\
&= \frac{1-\gamma^K}{K(1-\gamma)} \mathbb{E}_{\tau_m \sim \hat{\pi}_s} \left\{ \sum_{t \in T_m} \gamma^{t/K\alpha} A^{\tilde{\pi}_m}(\mathbf{o}_t^{m \rightarrow s}, \mathbf{g}_t^{m \rightarrow s}) \right\} + \eta(\hat{\pi}_s) = \\
&= \left(\frac{1-\gamma^K}{K(1-\gamma)} + \frac{1}{k_m} \right) \mathbb{E}_{\tau_m \sim \hat{\pi}_s} \left[\sum_{t \in T_m} \gamma^{t/K\alpha} A^{\tilde{\pi}_m}(\mathbf{o}_t^{m \rightarrow s}, \mathbf{g}_t^{m \rightarrow s}) \right] + \eta(\tilde{\pi}) = k_s \eta_{adv}(\hat{\pi}_s) + \eta(\tilde{\pi})
\end{aligned} \tag{15}$$

where used [Lemma A.1](#) and the approximation holds when $\gamma_s \simeq \gamma \simeq 1$ and K is not too large. Following [Corollary A.2](#), since $k_s > 0$, maximizing the objective of each sub-manager $\eta_s(\pi_s)$ implies maximizing the global objective $\eta(\hat{\pi}_s)$. This proves [Theorem 3.2](#). \square

A.4 Proof of Theorem 3.3

In general, goals received by workers are fixed for α steps. Therefore, the reward (i.e., the advantage function of the sub-manager) in Eq. 5 is constant in the α steps where the worker acts. We provide separate proofs of Theorem 3.3 for static and dynamic hierarchies to ease the presentation, despite the final result being the same.

Proof - Static Hierarchies. Let us denote as s the supervisor of the worker w . The objective for w reads:

$$\begin{aligned}
\eta_w(\pi_w) &= \mathbb{E}_{\mathbf{S}_0 \mathbf{o}_0^w} \mathbb{E} [V^{\hat{\pi}_w}(\mathbf{o}_0^w)] = \mathbb{E}_{\tau_w \sim \hat{\pi}_w} \left\{ \sum_{t=0}^{\infty} \gamma_w^t r_t^w \right\} = \mathbb{E}_{\tau \sim \hat{\pi}_w} \left\{ \sum_{t=0}^{\infty} \gamma_w^t \frac{A^{\hat{\pi}_s}(\mathbf{o}_t^{s \rightarrow w}, \mathbf{g}_t^{s \rightarrow w})}{\alpha} \right\} = \\
&= \frac{1}{\alpha} \mathbb{E}_{\tau \sim \hat{\pi}_w} \left\{ \sum_{t \in T_s} \gamma_w^t \sum_{i=0}^{\alpha-1} \gamma_w^i \left[r_t^{s \rightarrow w} + \gamma_s V^{\hat{\pi}_s}(\mathbf{o}_{t+\alpha}^{s \rightarrow w}) - V^{\hat{\pi}_s}(\mathbf{o}_t^{s \rightarrow w}) \right] \right\} \simeq \\
&\simeq \frac{1 - \gamma^\alpha}{\alpha(1 - \gamma)} \mathbb{E}_{\tau \sim \hat{\pi}_w} \left\{ \sum_{t \in T_s} \gamma^{t/\alpha} r_t^{s \rightarrow w} - V^{\hat{\pi}_s}(\mathbf{o}_0^{s \rightarrow w}) \right\} = \\
&= \frac{1 - \gamma^\alpha}{\alpha(1 - \gamma)} \left\{ \mathbb{E}_{\tau \sim \hat{\pi}_w} \left[\sum_{t \in T_s} \gamma^{t/\alpha} r_t^{s \rightarrow w} \right] - \mathbb{E}_{\mathbf{S}_0 \mathbf{o}_0^{s \rightarrow w}} \left[V^{\hat{\pi}_s}(\mathbf{o}_0^{s \rightarrow w}) \right] \right\} = \\
&= \frac{1 - \gamma^\alpha}{\alpha(1 - \gamma)} \left[k_s \eta_{adv}(\hat{\pi}_w) + \eta(\hat{\pi}) - \eta_s(\hat{\pi}_s) \right] = k_w \left[k_s \eta_{adv}(\hat{\pi}_w) + \eta(\hat{\pi}) - \eta_s(\hat{\pi}_s) \right] \quad (16)
\end{aligned}$$

where we used Theorem 3.2 and assumed $\gamma_w \simeq \gamma_s \simeq 1$ and α not too large. Since $k_w > 0$ and given Corollary A.2, maximizing the objective $\eta_w(\pi_w)$ of each worker implies maximizing the global objective $\eta(\hat{\pi}_w)$. This proves Theorem 3.3 in the case of a static hierarchy. \square

Proof - Dynamic Hierarchies. As we did for (sub-)managers, we evaluate the trajectory by keeping the hierarchy fixed. In other words, we consider the sequence of rewards associated with the sub-manager s_0 at the first step $t = 0$. Nevertheless, defining the future rewards when using a dynamic hierarchy is an ill-posed problem. Indeed, observations $\mathbf{o}_t^{s_0 \rightarrow w}$ at future steps $t > 0$ are not properly defined when the actual supervisor δ_t of w at time t is not s_0 . However, as reported at the end of Sec. 3.3, using a 3-level hierarchy restricts the model's applicability to cooperative problems. We then expect sub-managers to be cooperative too, and we can make them benefit from the value functions of future supervisors $\delta_{t>0}$ of w . Therefore, fixed the initial sub-manager s_0 and given a future step t , we can define the advantage-like reward of w under the supervision of the sub-manager s_0 as:

$$r_t^w = \frac{1}{\alpha} \left[r_t^{s_0 \rightarrow w} + \gamma_s V^{\hat{\pi}_s}(\mathbf{o}_{t+\alpha}^{\delta_{t+\alpha} \rightarrow w}) \llbracket t + \alpha \leq t^* \rrbracket - V^{\hat{\pi}_s}(\mathbf{o}_t^{\delta_t \rightarrow w}) \llbracket t \leq t^* \rrbracket \right] \quad (17)$$

where the Iverson brackets $\llbracket \cdot \rrbracket$ allow for truncating the influence of the goal $\mathbf{g}_0^{s_0 \rightarrow w}$ to the sequence of future sub-managers after the truncation step t^* . The objective for w reads:

$$\begin{aligned}
\eta_w(\pi_w) &= \mathbb{E}_{\mathbf{S}_0 \mathbf{o}_0^w} \mathbb{E} [V^{\hat{\pi}_w}(\mathbf{o}_0^w)] = \mathbb{E}_{\tau_w \sim \hat{\pi}_w} \left\{ \sum_{t=0}^{\infty} \gamma_w^t r_t^w \right\} = \\
&= \frac{1}{\alpha} \mathbb{E}_{\tau \sim \hat{\pi}_w} \left\{ \sum_{t \in T_s} \gamma_w^t \sum_{i=0}^{\alpha-1} \gamma_w^i \left[r_t^{s_0 \rightarrow w} + \gamma_s V^{\hat{\pi}_s}(\mathbf{o}_{t+\alpha}^{\delta_{t+\alpha} \rightarrow w}) \llbracket t + \alpha \leq t^* \rrbracket - V^{\hat{\pi}_s}(\mathbf{o}_t^{\delta_t \rightarrow w}) \llbracket t \leq t^* \rrbracket \right] \right\} \simeq \\
&\simeq \frac{1 - \gamma^\alpha}{\alpha(1 - \gamma)} \mathbb{E}_{\tau \sim \hat{\pi}_w} \left\{ \sum_{t \in T_s} \gamma^{t/\alpha} r_t^{s_0 \rightarrow w} - V^{\hat{\pi}_s}(\mathbf{o}_0^{s_0 \rightarrow w}) \right\} = \\
&= \frac{1 - \gamma^\alpha}{\alpha(1 - \gamma)} \left\{ \mathbb{E}_{\tau \sim \hat{\pi}_w} \left[\sum_{t \in T_s} \gamma^{t/\alpha} r_t^{s_0 \rightarrow w} \right] - \mathbb{E}_{\mathbf{S}_0 \mathbf{o}_0^{s_0 \rightarrow w}} \left[V^{\hat{\pi}_s}(\mathbf{o}_0^{s_0 \rightarrow w}) \right] \right\} = \\
&= \frac{1 - \gamma^\alpha}{\alpha(1 - \gamma)} \left[k_s \eta_{adv}(\hat{\pi}_w) + \eta(\hat{\pi}) - \eta_s(\hat{\pi}_s) \right] = k_w \left[k_s \eta_{adv}(\hat{\pi}_w) + \eta(\hat{\pi}) - \eta_s(\hat{\pi}_s) \right] \quad (18)
\end{aligned}$$

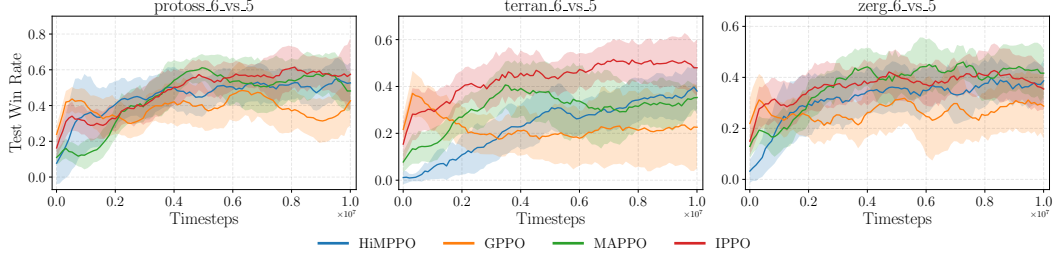


Figure 7: Results of the models for three maps of the *SMACv2* environment. We report the average test win rate and sample standard deviation of 4 runs for each configuration.

where we used [Theorem 3.2](#) and assumed $\gamma_w \simeq \gamma_s \simeq 1$ and α not too large. Since $k_w > 0$ and given [Corollary A.2](#), maximizing the objective $\eta_w(\pi_w)$ of each worker implies maximizing the global objective $\eta(\hat{\pi}_w)$. This proves [Theorem 3.3](#) in the case of a dynamic hierarchy. \square

Notice that a similar result for dynamic hierarchies can be achieved using an advantage-like reward defined in the time scale of the workers:

$$r_t^w = \frac{1}{\alpha} \left[r_t^{s_0 \rightarrow w} + \gamma_s V^{\pi_s}(\mathbf{o}_{t+1}^{\delta_{t+1} \rightarrow w}) \mathbb{I}[t+1 \leq t^*] - V^{\pi_s}(\mathbf{o}_t^{\delta_t \rightarrow w}) \mathbb{I}[t \leq t^*] \right] \quad (19)$$

where:

$$\mathbf{o}_{t+1}^{\delta_{t+1} \rightarrow w} = \begin{cases} \mathbf{o}_{t+\alpha}^{\delta_{t+\alpha} \rightarrow w} & \text{if } (t+1) \bmod \alpha = 0 \\ \mathbf{o}_t^{\delta_t \rightarrow w} & \text{otherwise} \end{cases} \quad (20)$$

By doing so, workers undergo the possible change of supervisor only at the last of the α steps. We conduct an empirical analysis of this aspect in [App. B.2](#).

B Additional Results

B.1 SMACv2

We consider 3 representative maps from *SMACv2*, namely *Protoss*, *Terran*, and *Zerg*. Since hyperparameters were not tuned for this environment, we adopted the simplified *6_vs_5* scenario for all the maps. For HiMPPO, the hierarchical scheme is a 2-level hierarchy where goals are propagated every 5 steps.

The results reported in [Fig. 7](#) show that graph-based representations do not improve performance in this environment, as GPPO performs poorly on average. Indeed, allies' features are already encoded in the observation space of each agent. Nevertheless, the 2-level hierarchical coordination mechanism allows HiMPPO to achieve performance comparable with IPPO and MAPPO despite being less sample efficient. Indeed, independent learning can achieve remarkable performance in similar environments [\[19\]](#), highlighting that additional information processing mechanisms might be redundant to solve these tasks. Since *SMACv2* is a fully cooperative environment, agents receive the same reward signal at each step. This makes it more challenging for our model to learn an effective hierarchical coordination mechanism, as goals received by the workers at the same step are associated with the same long-term return. Therefore, the manager has to learn to generate individual goals based on an implicit estimate of the local contributions of each worker. Nevertheless, despite the complex reward-assignment mechanism, HiMPPO achieves a competitive win rate percentage.

B.2 Analysis of the Truncation Value

As reported in [App. A.4](#), dynamic 3-level hierarchies require defining future value functions when considering the sequence of rewards associated with the sub-manager at the first step. In particular, according to the desired temporal resolution, they can be defined w.r.t. the time scale of

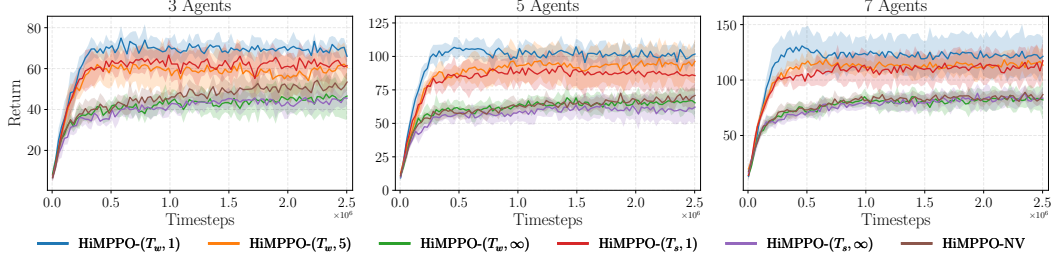


Figure 8: Results of the ablation study for different truncation schemes in the *VMAS Sampling* environment. For each configuration, we report the average return and sample standard deviation of 6 runs.

sub-managers (Eq. 17) or workers (Eq. 19). Moreover, future values $V\pi_s(o_t^{\delta_t \rightarrow w})$ can be set to zero after a certain truncation step t^* . This threshold accounts for the degree of cooperation of the sub-managers: low values of t^* imply that the sub-manager sending the goal relies only on first estimates of the value functions, and vice versa.

To investigate these aspects, we consider both the value-assignment methods as well as different truncation values. We denote as $\text{HiMPPO-}(T_s, t^*)$ and $\text{HiMPPO-}(T_w, t^*)$ the models where advantage-like rewards are defined using Eq. 17 and 19, respectively, with truncation step t^* . Furthermore, No Values HiMPPO (HiMPPO-NV) denotes the advantage-like rewards without value functions, i.e., with truncation step $t^* = 0$. The results reported in Fig. 8 show that for low t^* the model performs better. However, completely truncating the value functions from the workers' reward (HiMPPO-NV) leads to lower returns. A possible explanation could be that the sub-managers' value functions are crucial for correctly guiding the workers toward the global objective. Nevertheless, high values of t^* make credit assignment more challenging, as it considers a longer sequence of value functions, which might be related to other sub-managers.

C Experimental Setting

C.1 Environments

C.1.1 Level-Based Foraging with Survival (LBFwS)

Starting from the codebase provided by Azmani et al. [7], we modify the Level-Based Foraging [3, 59] environment by further exacerbating the cooperative-competitive behavior of the agents to make it more challenging. In particular, once the food is within its range, an agent can eat it or deliver it to a landmark. Following the original version, the first strategy leads to an immediate, individual reward. Conversely, delivering food extends the episode duration: even if it gives no individual reward in the short term, it is beneficial for the team of agents, allowing them to have more time to collect food and achieve a potentially higher cumulative return. This behavior mirrors community food-sharing dynamics, where individuals contribute to the group's survival. This extended version of the environment, named *Level-Based Foraging with Survival* (LBFwS), aggravates the cooperative-competitive dynamics: in addition to evenly dividing the reward when cooperating [7], trying to deliver food can be harmful to the agent, as other agents might act greedily to maximize their individual reward instead of adopting the group's survival strategy. Furthermore, delivering food requires a complex and temporally extended series of actions, i.e., collecting the item, transporting it to the landmark, and releasing it, making the problem challenging in terms of exploration and credit assignment. We now illustrate the differences of the LBFwS environment with respect to the version of Azmani et al. [7].

Observation Space Compared with the original setting, each agent has an additional grid representing its local observation of the landmark layer, i.e., if the landmark is within its sight range and if the

agent (or its neighbors) is carrying food. Furthermore, each i -th agent observation is augmented by a 4-dimensional vector $(t, t_s, x_{lm}^i, y_{lm}^i)$, where t is the time passed so far (normalized with the absolute maximum number of steps T), t_s is the survival time (which can be incremented by delivering food) normalized with the survival initial value $T_s < T$, and x_{lm}^i and y_{lm}^i are the normalized relative displacements in the (x, y) axis w.r.t. the landmark position. Therefore, agents are always aware of the position of the landmark, which is static and placed in the middle of the map.

Action Space We add two possible actions to the original 6-dimensional discrete action space, i.e., PICK and DELIVERY. Agents that take the EAT action while carrying food consume the item. Conversely, if an agent attempts to pick up food while already carrying an item, that action is automatically canceled. In general, unsuccessful pickup or delivery attempts do not affect the current state of the environment. We remark that items of higher value need the cooperation of multiple agents to be consumed.

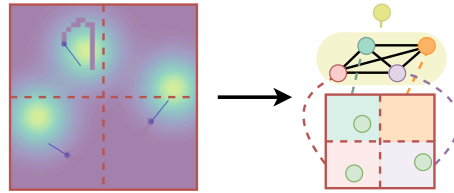
Reward We highlight that the reward structure does not change, as delivering food does not affect the agents' rewards.

Transition Dynamics The survival and absolute counters start at $t_s = T_s$ and $t = 0$, respectively. When an agent successfully delivers food to the landmark, the counter t_s is increased by the value of the food, incremented by a factor k_{surv} . At each step, the survival (absolute) counter decreases (increases) by one, and the episode terminates either when $t_s = 0$ or $t = T$. Each agent can carry only one item, while unsuccessful deliveries do not imply negative rewards or changes in the item that is carried. Furthermore, an agent can decide to eat the food that is being carried.

Experimental Setup In our experiments, we fix the maximum values for survival and absolute steps as $T_s = 100$ and $T = 500$, respectively. Furthermore, we set the constant for the survival time increment to $k_{surv} = 10$: as an example, successfully delivering an item with value 4 leads to an increment of t_s of 40 time steps. In our experiments, we consider three different maps of different sizes and amounts of resources. In particular, we consider: 1) *LBFwS-Easy*, where the map is a 9×9 grid and there are 4 food items of level 1 and 4 food items of level 2; 2) *LBFwS-Medium*, where the map is a 12×12 grid and there are 5 food items of level 1 and 5 food items of level 2; 3) *LBFwS-Hard*, where the map is a 15×15 grid and there are 6 food items of level 1 and 6 food items of level 2. We remark that level 2 items need the cooperation of at least 2 agents to be picked up. For all the scenarios, the number of agents is set to 10. Despite having more food available when the difficulty increases, the grid size makes exploration and temporal credit assignment more challenging. Furthermore, this difficulty is compounded by the sparser distribution of food items in the space, which makes successful deliveries to the landmark more complex.

C.1.2 Hierarchical Graph Generation in Sampling

As aforementioned, in the *VMAS Sampling* scenario, we consider a 3-level dynamic graph \mathcal{G}_t^* for the HiMPPO model. This graph is built according to the current state \mathbf{S}_t of the environment. In particular, the grid is divided into 4 quadrants, and each worker (agent) is assigned to a single sub-manager based on its position. Refer to Fig. 9 for a visual example.



We augment the initial representation $\mathbf{h}_t^{s,0}$ of each sub-manager s with a one-hot vector to give an explicit indexing to the sub-managers.

However, we remark that this does not prevent the transferability of the architecture to systems with a different number of agents. There might be states where a sub-manager has no workers underneath, e.g., the up-right (orange) sub-manager of Fig. 9. In such cases, the top-level manager sends a goal to account for possible future changes in the dynamics, but the corresponding trajectory is discarded in the learning procedure.

Figure 9: Building the 3-level hierarchical graph \mathcal{G}_t^* at time t for the *VMAS Sampling* scenario.

C.2 Hardware and Software

The code for the model and the baselines was developed in *Python 3.10* and made use of *NumPy* [30], *PyTorch* [60], and *PyTorch Geometric* [23]. We implemented the PPO-based models starting from a public repository [9] (MIT License). For the environments, we used the official repositories [12, 7, 22] and the utility functions of *TorchRL* [15] (MIT License). *VMAS Sampling* uses a GPL3.0 License, while *SMACv2* uses MIT License. Experiment configurations were managed using *Hydra* [71] and tracked using *Neptune.ai*¹. Simulations were run on two workstations equipped with AMD EPYC 7513 and Intel Xeon E5-2650 CPUs. Depending on the model, training times take 2 – 10 hours for *VMAS Sampling*, 8 – 15 hours for *LBFwS*, and 1 – 20 days for *SMACv2*.

C.3 Implementation Details

To achieve the optimal configurations, for all the models we performed a grid search on *LBFwS-Medium* environment, focusing mainly on latent dimensions, hidden units, and learning rates. For vector-based observations, encoders are designed as a linear layer. For grid-based observations (*LBFwS*), we use a Convolutional Network, where we tuned the number of output features and kernel size; then, two linear encoders process the flattened space and time observations to obtain the desired latent dimension. All the models use a discount factor $\gamma = 0.99$ and the Adam [39] optimizer. For PPO-based models, the entropy coefficient and clipping values are fixed to 0.01 and 0.2, respectively. Advantages are approximated using Generalized Advantage Estimator [63] $GAE(\gamma, \lambda)$: λ is set to 0.95, except for upper levels of HiMPPPO (i.e., (sub-)managers), that is fixed to 0. For the *VMAS Sampling* environment, we tuned the initial action standard deviation, which then decays by 0.05 every $2.5 \cdot 10^5$ steps up to 0.1. Policies are updated for 30 epochs with a batch size of 128 every 40 (*LBFwS* and *SMACv2*) and 32 (*VMAS Sampling*) episodes. Unless constrained by the action space, we use ReLU as the activation function. For HiMPPPO and GPPO, the agents' message function is implemented as an MLP with 64 hidden units in *LBFwS* and *SMACv2*, and as a Graph Convolutional Network (GCN) [40] in *VMAS Sampling*. When employing a 3-level hierarchy (*VMAS Sampling*), HiMPPPO uses an MLP with 64 hidden units as the message function for sub-managers.

HiMPPPO All the hyperparameters are shared between the different levels of the hierarchy, except for λ . We used a learning rate of 0.0001 and 0.0005 for actor and critic, respectively. To keep the number of parameters bounded, all the functions are implemented as linear layers, except for the message function that has 64 hidden units; we consider 1 round of message passing. Embedding and goal dimensions are fixed to 64. The goal is sampled from a Gaussian policy with an initial standard deviation of 0.5. For the *VMAS Sampling* environment, where continuous actions are required, we use the same value for the initial standard deviation. For *LBFwS*, we consider a 2-dimensional convolution with 8 output channels and kernel size 2. To ensure stability among different hierarchy levels, we clip gradients to the range $[-1, 1]$.

GPPO This model does not share modules between the actor and critic. The actor and critic have a hidden dimension of 64 and 32, respectively. We used a learning rate of 0.0001 for both. The embedding dimension is fixed to 64. We perform 2 message-passing rounds. For the *VMAS Sampling* environment, the initial standard deviation is set to 0.5. For *LBFwS*, we consider a 2-dimensional convolution with 32 output channels and kernel size 2. Gradients are clipped to the range $[-5, 5]$.

MAPPO In this baseline, the embedding dimension is 128. The actor has 128 hidden units, whereas the critic is implemented as a linear layer. We used a learning rate of 0.0005 and 0.0001 for actor and critic, respectively. For the *VMAS Sampling* environment, the initial standard deviation is set to 0.5. For *LBFwS*, we consider a 2-dimensional convolution with 8 output channels and kernel size 2. Gradients are clipped to the range $[-5, 5]$.

IPPO This model uses an embedding dimension of 64. The actor and critic have 128 and 64 hidden units, whereas the learning rates have values 0.0005 and 0.00001, respectively. For the *VMAS Sampling* environment, the initial standard deviation is set to 0.8. For *LBFwS*, we consider a 2-dimensional convolution with 8 output channels and kernel size 2. Gradients are clipped to the range $[-5, 5]$.

¹<https://neptune.ai/>