

# Heterogeneous Value Decomposition Policy Fusion for Multi-Agent Cooperation

Siyang Wang<sup>1</sup>, Yang Zhou<sup>2</sup>, Zhitong Zhao<sup>2</sup>, Ruoning Zhang<sup>2</sup>,  
Jinliang Shao<sup>1</sup>, Wenyu Chen<sup>2</sup>, Yuhua Cheng<sup>1\*</sup>

<sup>1</sup>School of Automation Engineering, University of Electronic Science and Technology of China

<sup>2</sup>School of Computer Science and Engineering,

University of Electronic Science and Technology of China

siyangwang@uestc.edu.cn, {zhouy, zhaozhitong, zhangruoning}@std.uestc.edu.cn,  
{jinliangshao, cw, yhcheng}@uestc.edu.cn,

## Abstract

Value decomposition (VD) has become one of the most prominent solutions in cooperative multi-agent reinforcement learning. Most existing methods generally explore how to factorize the joint value and minimize the discrepancies between agent observations and characteristics of environmental states. However, direct decomposition may result in limited representation or difficulty in optimization. Orthogonal to designing a new factorization scheme, in this paper, we propose Heterogeneous Policy Fusion (HPF) to integrate the strengths of various VD methods. We construct a composite policy set to select policies for interaction adaptively. Specifically, this adaptive mechanism allows agents' trajectories to benefit from diverse policy transitions while incorporating the advantages of each factorization method. Additionally, HPF introduces a constraint between these heterogeneous policies to rectify the misleading update caused by the unexpected exploratory or suboptimal non-cooperation. Experimental results on cooperative tasks show HPF's superior performance over multiple baselines, proving its effectiveness and ease of implementation.

## 1 Introduction

Multi-Agent Reinforcement Learning (MARL) [Jaques *et al.*, 2019; Mei *et al.*, 2023], derived from reinforcement learning methods [Sutton and Barto, 2018], has proven to be successful in various tasks that involve cooperation and collaboration among multiple agents, such as coordination of robots [Perrusquía *et al.*, 2021], autonomous Unmanned Aerial Vehicles (UAVs) [Zhang *et al.*, 2023], and traffic signal control [Yang *et al.*, 2023]. These tasks often require agents to possess the ability to make distributed sequential decisions while considering collaboration with each other to collectively achieve a common goal. However, learning such effective cooperation policies among these distributed agents remains a major challenge due to the partial observability and non-stationarity in the environment result-

ing from the changing policies [Oliehoek and Amato, 2016; Tan, 1993].

A fully centralized information processing agent seems to be a good solution to address the above issues. Nevertheless, it encounters the challenge of an exponentially expanding joint action space with the increase in the number of collaborative agents. To cope with this situation, the Centralized Training with Decentralized Execution (CTDE) framework was put forward. The MARL algorithms embrace CTDE to train the agents' policies with access to the global state, mitigating the non-stationarity and unstable training. Meanwhile, decentralized execution enables scalability by depending on local information. Within this framework, value-decomposition (VD) algorithms estimate the true joint action value function by factoring it into decentralized utility functions and satisfying the Individual-Global-Max (IGM) principle. Additionally, IGM ensures consistency between optimal local and global actions. The core idea is that if the factored versions are accurate, the optimal joint actions can be determined through local greedy operations. However, a particular type of VD approach faces representation limitations by restricting its mixing network parameters to be positive to meet the IGM condition, which simplifies learning but can hinder finding the optimal joint actions.

Another type of VD method endeavors to construct a surrogate target by leveraging the local optimal actions of distributed agents as prior conditions, and narrows the disparity between the surrogate objective and the true joint action-value function to overcome representation limitations. However, in more complex cooperative tasks, this gap may never be fully bridged. It is challenging for the model to simultaneously learn and estimate the true and surrogate joint action-value functions during dynamic interactions between agents, requiring careful calibration of the update weights for both objectives, or incorporating additional external features to enhance the accuracy of joint value function estimation. These limitations result in VD methods of this type either requiring more complex designs or encountering problems like low training efficiency and poor agent collaboration.

To overcome the aforementioned limitations, this paper proposes a Heterogeneous Policy Fusion (HPF) scheme to integrate the benefits of different types of VD methods. Our work shares the same objective of effectively enhancing the

\*Corresponding author

training efficiency of MARL algorithms. However, we focus on constructing a novel composite policy to interact with the environment and gather more valuable data instead of developing a new VD scheme. Specifically, HPF expands different VD methods into a policy set with an adaptive composition based on the value estimation, which determines the degree of participation of each candidate policy. To prevent unexpected exploration and suboptimal non-cooperation, an instructive constraint based on KL divergence is imposed between the utility functions of the candidate VD methods. This constraint prevents excessive divergence between their policies, helping to correct the misleading updates in the model. Our contributions can be summed up as threefold:

- we analyze the connections and differences among different types of VD methods, and highlight the value of properly integrating these methods in order to enjoy the best of each part, a perspective that is unique and orthogonal to developing a completely new VD scheme.
- we propose a simple VD policy extension scheme to leverage the benefits of heterogeneous VD policies, and show that HPF can effectively enhance the training efficiency of the collaboration model.
- we verify the effectiveness of HPF by comparing it against different VD methods across several cooperative multi-agent tasks, to demonstrate its superiority to tackle issues like *Relative Overgeneralization*, showcasing its remarkable performance.

## 2 Related Works

Value decomposition is effective in addressing challenges such as partial observability, algorithmic instability, and credit assignment in cooperative environments. Our focus lies in integrating the advantages of VD methods based on the Individual-Global-Maximum (IGM) [Son *et al.*, 2019] principle. Categorized by the different ways of satisfying IGM, the methods fall into two groups: one constrained by network structure, the other constructing a surrogate target and involving information replenishment.

The first category of VD methods decomposes the joint action-value function by imposing constraints on the network structure, such as *additivity* or *monotonicity*. Following the IGM principle, VDN [Sunehag *et al.*, 2017] decomposes the joint value function  $Q_{tot}$  into the sum of individual value functions  $\sum [Q_i]_{i=1}^n$ , neglecting additional information available during training. While QMIX [Rashid *et al.*, 2018] imposes monotonicity constraints on the mixing network, making the  $Q_{tot}$  and  $Q_i$  satisfy the IGM principle. Qatten [Yang *et al.*, 2020] theoretically derives the process of decomposing the joint Q-value ( $Q_{tot}$ ) into local Q-values ( $Q_i$ ). VGN [Wei *et al.*, 2022] takes into account the influence of agents with the minimum Q-value and utilizes graph attention networks to introduce dynamic relationships between agents. GraphMIX [Naderializadeh *et al.*, 2020] introduces graph network and attention mechanisms to the QMIX, integrating the information of agents through the edge weights controlled by attentional relationships between agents. TransMix [Khan *et al.*, 2022] propose a Transformer-based mixture network, emphasizing the extraction of global and local

contextual interactions among  $Q_i$ , historical trajectories, and global states. The recent MDQ [Ding *et al.*, 2023a] addresses cooperative, competitive, or mixed tasks by combining mean-field theory with VD. Inspired by the human nervous system, HAVEN [Xu *et al.*, 2023] devised a two-level QMIX strategy for inter-level and intra-level coordination. Restricting the network parameters to achieve IGM can lead to representation limitation of the joint action-value function, making it hard to cover the optimal joint actions.

Another category of VD methods satisfies IGM by constructing a surrogate target and supplementing additional information to approximate the true joint action-value function. QTRAN [Son *et al.*, 2019] introduces a soft regularization constraint, wherein the theoretically learned overall reward function  $Q_{tran}$  is required to be equal to the true value of the overall reward function  $Q_{tot}$  only when the optimal action  $\bar{u}$  is taken. While QTRAN++ [Son *et al.*, 2020] explicitly specifies the relationship  $Q_{tot}(s, \tau, \bar{u}) = Q_{tran}(s, \tau, \bar{u}) > Q_{tran}(s, \tau, u) > Q_{tot}(s, \tau, u)$  and introduces the corresponding loss for  $Q_{tot}$  to expedite training. WQMIX [Rashid *et al.*, 2020] uses weighted projection to improve QMIX and reduce the learning weight of non-optimal joint actions. QPLEX [Wang *et al.*, 2021] converts the IGM principle into consistency constraints for the advantage function, and uses a dual-competitive architecture to factorize  $Q_{tot}$ . ResQ [Shen *et al.*, 2022] uses the idea of residual decomposition to decompose the joint Q-value into the sum of a main function and a residual function, to solve the problem of limited representation. RQN [Pina *et al.*, 2022] learns an individualized factor for each agent signifying the relative importance of a specific Q-value trajectory. LOMAQ [Zohar *et al.*, 2022] proposes a scalable VD method to improve credit assignment by utilizing the learned local agent rewards. CIA [Liu *et al.*, 2023] introduces a new contrastive identity-aware learning method to explicitly encourage differentiation between credit levels of agents. MARGIN [Ding *et al.*, 2023b] decomposes the global mutual information into a weighted sum of local mutual information, making it seamlessly integrated with various VD approaches.

Although these methods can obtain more information from the environment or other agents, they usually require a carefully designed structure. Orthogonal to developing a new VD method, our proposed HPF amalgamates the strengths of the aforementioned VD approaches, aiming to achieve a highly sample-efficient and expressive VD approach.

## 3 Preliminaries

### 3.1 Dec-POMDPs

The formalism known as *Decentralized Partially Observable Markov Decision Processes* (Dec-POMDPs) [Oliehoek and Amato, 2016], characterized by the tuple  $G = (S, U, P, r, Z, O, n, \gamma)$ , is widely employed in delineating the dynamics of complete cooperative multi-agent endeavors. The system comprises  $n$  agents, denoted as  $n \in N \equiv 1, \dots, n$ , and the genuine environmental state  $s \in S$  encapsulates both the agents' collective knowledge and additional contextual features. At every timestep  $t$ , the environment transitions based on the state transition function

$P(s' | s, \mathbf{u}) : S \times \mathbf{U} \times S \rightarrow [0, 1]$ , where  $\mathbf{u} \in \mathbf{U} \equiv U^n$  encompasses the collectively generated actions with an action  $u_i \in U$  chosen by each agent. A global reward function  $r(s, \mathbf{u}) : S \times \mathbf{U} \rightarrow \mathbf{R}$  is shared among all agents, and each agent receives a partial observation  $z^i \in Z$  based on the observation function  $O(z^i | s, u^i) : S \times U \rightarrow Z$ . The primary goal for all agents is to jointly optimize the team reward  $R_t = \sum_{i=0}^T \gamma^i r_{t+i}$ , elucidated through the joint value function  $Q^\pi(s_t, \mathbf{u}_t) = \mathbb{E}_{s_{t+1:\infty}, \mathbf{u}_{t+1:\infty}} [R_t | s_t, \mathbf{u}_t]$ .

### 3.2 Value Function Decomposition

#### Implementation by network constraint

The connotation of value function decomposition is to establish a direct training connection between the joint state-action value function and the utility functions of the individual agents during centralized training. Thus, for a joint state-action value function  $Q_{tot}$ , there exist utility functions  $[Q_i]_{i=1}^n, i \in N$  for each of the agents that satisfy:

$$\begin{aligned} \arg \max_{\mathbf{u}} Q_{tot}(\tau, \mathbf{u}) = \\ (\arg \max_{u^1} Q_1(\tau^1, u^1) \dots \arg \max_{u^n} Q_n(\tau^n, u^n)). \end{aligned} \quad (1)$$

This equation can be described as  $[Q_i]_{i=1}^n$  satisfying the IGM criterion for the hypothetical decomposition-available joint-valued function  $Q_{tot}$ . Therefore, it is also asserted that  $Q_{tot}(\tau, \mathbf{u})$  can be decomposed into  $[Q_i(\tau_i, u_i)]_{i=1}^n$  or  $[Q_i]_{i=1}^n$  is a factor of the decomposition of  $Q_{tot}$ . Two most typical value function decomposition methods are VDN [Sunehag *et al.*, 2017] and QMIX [Rashid *et al.*, 2018], which satisfy the sufficient conditions as  $Q_{tot}^{VDN}(\tau, \mathbf{u}) = \sum_{i=1}^n Q_i(\tau^i, u^i)$  and  $Q_{tot}^{QMIX}(\tau, \mathbf{u}) = \sum_{i=1}^n |w_i| Q_i(\tau^i, u^i)$ . These two methods calculate the global action value function by introducing the mixing network as  $Q_{tot} = g(Q_1, Q_2, \dots, Q_n)$ , or factorize the centralized value function with the assurance of monotonicity property, with parameters restricted to be positive. However, the monotonicity constraint limits the representational capacity of the factored global action-value function, which may lead to the agents failing to find the correct optimal joint actions.

#### Implementation by constructing surrogate target

Another type of VD method requires constructing a surrogate objective to approximate the true joint action-value function, ensuring they are equal when taking the optimal joint action. Here, we denote the true joint action-value function as  $Q_{jt}$ , and the surrogate objective is estimated by the aforementioned VD methods with network constraints as  $Q_{tot}$ . QMIX [Rashid *et al.*, 2020] tries to pay more attention to the optimal joint action during the training process, and designs to update the true joint action value function as  $Q_{jt} \leftarrow w(s, \mathbf{u}) Q_{jt}(s, \tau, \arg \max_{\mathbf{u}} Q_{tot}(\tau, \mathbf{u}, s))$  with weights  $w(s, \mathbf{u})$  determined by whether the current joint actions achieves the optimal one, which is depicted as:

$$w(s, \mathbf{u}) = \begin{cases} 1 & \mathbf{u} = \mathbf{u}^* = \arg \max_{\mathbf{u}} Q_{jt}(s, \mathbf{u}) \\ & \text{or } Q_{tot}(s, \mathbf{u}) < Q_{jt}(s, \mathbf{u}), \\ \alpha & \text{otherwise.} \end{cases} \quad (2)$$

While QPLEX designs an Advantage-based IGM principle, which loosens the expressiveness of the IGM-class factored global action value function through a duplex dueling architecture as:

$$\begin{aligned} Q_{tot}(\tau, \mathbf{u}) = V_{tot}(\tau) + A_{tot}(\tau, \mathbf{u}) \\ \text{with } V_{tot}(\tau) = \max_{\mathbf{u}'} Q_{tot}(\tau, \mathbf{u}'). \end{aligned} \quad (3)$$

## 4 Methodology

In this section, we will introduce the policy fusion pattern HPF in detail. Specifically, we begin by introducing the motivation behind our proposed scheme and then delve into the in-depth analysis of the composite policy in HPF. Additionally, we design an instructive constraint between the utility functions of composite policy, which aims to correct the learned factorized utilities and provide the VD learner with a more accurate estimation potential. Finally, we depict the overall implementation of HPF at the end of this section.

### 4.1 Motivation: connections and gaps between different VD policies

Our proposed fusion framework is founded on the VD paradigm. It decomposes the centralized value function into a combinatorial form of the utility functions for the collaborative model to learn. The aim of this factorization is to align the optimal actions of the centralized value function with the local optimal actions of individual agents. Then the collaborative model can update the utility functions by maximizing the  $Q$ -value while meeting the IGM criterion.

As summarized in Section 3, even though the methods of the two existing types of VD policies are both intended to more effectively align the optimal joint actions of the centralized value function and the local optimal actions of each agent, there are still obvious commonalities and differences among various approaches. These relationships and distinctions can be summed up as the following points.

#### For VD by network constraint

- *Direct centralized value function decomposition*: it requires fewer neural network components to be optimized, resulting in a faster optimization speed and higher efficiency.
- *Hard constraints on mixing network parameters*: the assurance of the IGM criterion is realized by limiting the parameters and structure of the mixing network.
- *Limited representation capacity*: the restrictions of the mixing network prevent the correct fitting of the negative correlation between the centralized value function and the utility functions.

#### For VD by constructing surrogate target

- *Surrogate target reservation*: the constructed surrogate target is incorporated as a new element into the original VD procedure, which may bring about relatively low training efficiency.
- *Soft constraints on optimization*: both the optimization and alignment process of actions are softly constrained

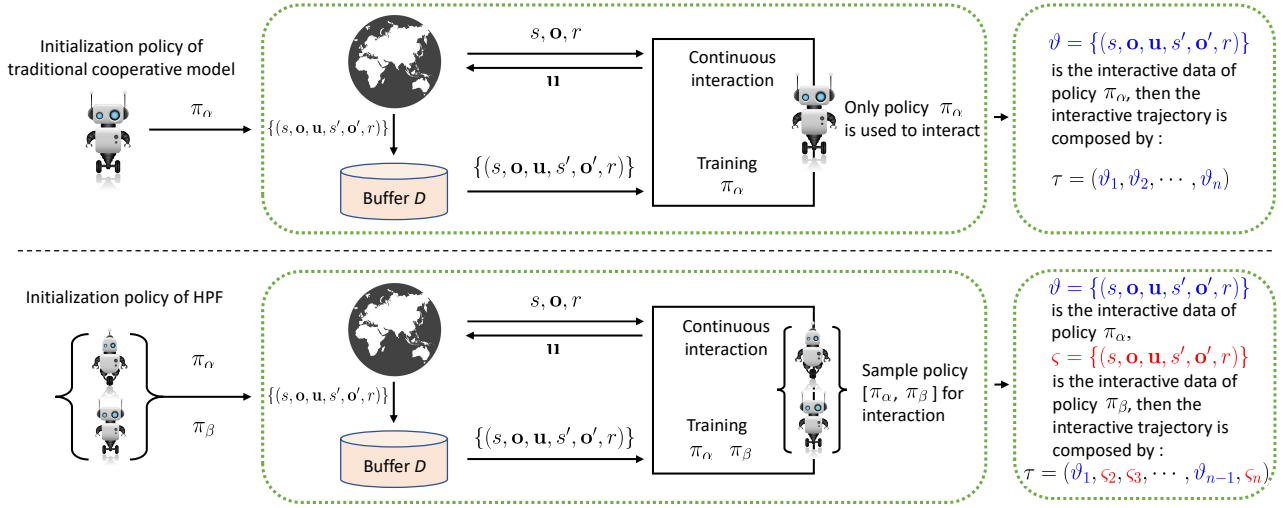


Figure 1: **The illustration of the distinction between HPF and traditional VD methods.** The traditional scheme directly aligns the optimal joint action and optimizes the central value function with the presupposed VD policy itself. The proposed HPF integrates the benefits of different types of VD policies, and expands them into a policy set to sample the experiences for capturing further performance improvement. These VD policies both participate in the interactions with the environment and learning in an adaptive manner.

by *MSE loss*, which may make it difficult to train the model with multiple objectives to be optimized.

- *Full representation capacity*: the surrogate target secludes the original true centralized value function objective, which maintains the full representation capacity of the neural network, enabling it to fully fit the complicated reward distribution.

Besides the optimization process and algorithmic connections, the two types of VD policies are also complementary to each other in terms of strengths and weaknesses. The VD policies with network constraints are more efficient since they have fewer components to optimize, but they may not be able to find the true optimal joint action with the limited representation. While another type of method enjoys the potential to find the true optimal joint action, but is comparatively much less sample efficient. This inevitably poses a question:

“Is there a way to integrate the benefits of the two types of VD policies with the interactions between the agents and environments?”

Because of the connections and complementary strengths above, instead of treating the different types of policies as two isolated aspects, it is more natural to connect both in pursuit of a performant policy in practice.

## 4.2 Policy Extension and Adaptive Composition

**Policy Extension.** The policy extension in HPF is direct and readily combined with existing VD methods, which is illustrated in Figure 1. Before the collaborative model interacts with the environment, we consider the representation-limited type of VD policy as a candidate policy. It is then extended into the surrogate-target type of policy, forming a candidate set. Subsequently, the collaborative model adaptively synthesizes a composite policy from this set to interact with the environment, which results in the mixed sample trajectory simultaneously containing interaction experiences derived from the

two types of VD approaches. Since the mixed sample trajectory potentially includes policy information from two different VD policies, both of them would preserve their respective characteristics while improving the shortcomings of each when they sample these experiences to train the collaborative model. The expanded policy set  $\Pi$  is depicted as:

$$\Pi = [\pi_\alpha, \pi_\beta], \quad (4)$$

where  $\pi_\alpha$  represents the VD policy with surrogate target, and  $\pi_\beta$  is VD policy with limited representation. In the context of the extended set of policies formed by existing VD policies, the model continues to meet the IGM criterion, which is shown in Proposition 1. The policies in  $\Pi = [\pi_\alpha, \pi_\beta]$  will all get involved in the interactions and learning in a collaborative and adaptive manner as detailed in the next part.

**Proposition 1.** *If  $\Pi = [\pi_\alpha, \pi_\beta]$  is an extended policy set formed by existing various VD policies, then  $\Pi$  still satisfies the IGM criterion.*

*Proof.* See Section A of Supplementary Materials.  $\square$

**Adaptive Composition.** Both VD policies in  $\Pi$  will form a single composite policy  $\tilde{\pi}$  from the heterogeneous VD policies in  $\Pi$ . More specifically, given the current state  $s$  and the joint observations  $\mathbf{o}$ , HPF first samples actions for each member of the policy set  $\Pi$  and then calculates the utility functions of each candidate VD policy. Since these utilities will contribute to determining the joint action-value functions, which is the lead of VD policy to update, we use them as a proxy to composite the final policy with probability based on them. For example, we estimate their values at the current state as  $\mathbb{Q} = \{Q^{\pi_\alpha}, Q^{\pi_\beta}\}$ . While each agent chooses the action with  $\epsilon$ -greedy as the interactive action, the set of utility functions is represented as  $Q^{\pi_\alpha} = \{Q_1^{\pi_\alpha}, Q_2^{\pi_\alpha}, \dots, Q_n^{\pi_\alpha}\}$  and  $Q^{\pi_\beta} = \{Q_1^{\pi_\beta}, Q_2^{\pi_\beta}, \dots, Q_n^{\pi_\beta}\}$ . Subsequently, HPF then

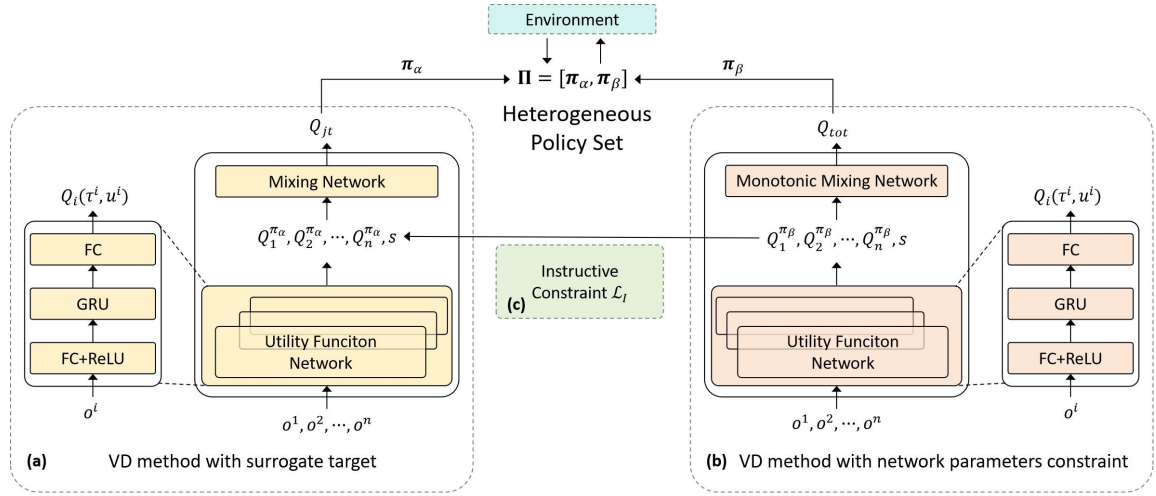


Figure 2: **The architecture of HPF.** (a) The VD method with surrogate target. (b) The VD method with network parameters constraint. (c) The instructive constraint between heterogeneous utility functions. The policies of both VD methods constitute a composite policy set and interact with the environment after sampling.

constructs a Boltzmann policy-based categorical distribution for selecting the final interactive policy:

$$P_{\mathbf{w}}[k] = \frac{\exp(Q^{\pi_k}/\eta)}{\sum_k \exp(Q^{\pi_k}/\eta)}, \quad \forall k \in [\alpha, \beta] \quad (5)$$

where  $\eta$  is the temperature of the Boltzmann distribution. Then we sample from  $\mathbf{w} \sim P_{\mathbf{w}}$  for integrating the composite policy to determine which policy will be implemented during unroll in interaction. The composite policy  $\tilde{\pi}$  can be conceptually represented in the following manner:

$$\tilde{\pi}(\mathbf{u} \mid \mathbf{o}) = [\delta_{\mathbf{u} \sim \pi_\alpha}, \delta_{\mathbf{u} \sim \pi_\beta}] \cdot \mathbf{w}, \quad \mathbf{w} \sim P_{\mathbf{w}} \quad (6)$$

where  $\mathbf{w}$  is a one-hot vector, indicating the policy that is selected for the joint observation  $\mathbf{o}$ .  $\delta_{\mathbf{u} \sim \pi}$  denotes the Dirac delta distribution centered at  $\mathbf{u}$  which is sampled from  $\pi$ .

The estimation method of  $Q^{\pi_k}$  in Equation (5) is crucial in the computation of  $P_{\mathbf{w}}$  because it influences the sampling process and fusion level of different heterogeneous VD policies within the  $\Pi$ . Therefore, we design two principles to estimate  $Q^{\pi_k}$ : (1) it uses the sum of utility functions in  $\Pi$  as the sampling criterion inspired by VDN [Sunhag *et al.*, 2017]; (2) it uses the estimation of joint optimal action value functions for the sampling, which takes into account the differences in the representational capabilities [Rashid *et al.*, 2020]. The detailed calculation is outlined as follows:

$$\text{(Additive)} \quad Q^{\pi_k} = \sum_{i=1}^n Q_i^{\pi_k} \quad (7)$$

$$\text{(Optimistic)} \quad Q^{\pi_k} = Q_{jt}^{\pi_k}(\tau, \hat{\mathbf{u}}^*, s),$$

where  $k \in [\alpha, \beta]$ ,  $i \in \{1, \dots, n\}$  indicate the agent ID, and  $\hat{\mathbf{u}}^* = \operatorname{argmax}_{\mathbf{u}} Q_{tot}(\tau, \mathbf{u}, s)$  means the set of maximal local actions. After that, the model interacts with the environment and forms a mixed trajectory to store in the replay buffer. During the training phase, the model samples a batch of sampled trajectories and trains both types of VD policies in

$\Pi$  simultaneously. Since the chosen interaction policy at each time step may differ, these trajectories may contain segments of experiences from different types of policies, intertwining the mixed experiences throughout the entire trajectory.

In general, HPF shares the same core as VD methods, aiming to enhance the training efficiency of collaborative policy. Orthogonal to designing a new factorization scheme, HPF focuses on integrating the advantages of existing VD methods using a policy fusion scheme. This also leverages the mixed experiences to guide and promote mutual interactions between the candidate VD policies in  $\Pi$ . In comparison to existing VD methods, HPF possesses several advantages:

- **Independence of candidate interaction policies:** The policies in the set  $\Pi = [\pi_\alpha, \pi_\beta]$  are mutually independent. Since the composite policy is implemented through sampling, the chosen interaction policy by the model remains unaffected by other policies in  $\Pi$ .
- **Flexibility in the form of interaction policies:** The policies in the set  $\Pi$  vary in type and network architecture. Since HPF aims to integrate the advantages of diverse policies, then  $\pi_\alpha$  and  $\pi_\beta$  are assigned two distinct types of policies within the VD framework. In a broader sense, HPF can flexibly adopt any method, allowing policies  $\pi_\alpha$  and  $\pi_\beta$  to differ in form, algorithm, or category based on the user's integration needs.
- **Adaptive policy selection:** Every policy within the set  $\Pi$  is available for interaction with the environment, engaging adaptively.

### 4.3 Instructive Constraint

Despite policy sampling, HPF's composite policy effectively leverages heterogeneous VD methods. However, VD methods in  $\Pi$  enforcing the IGM criterion via parameter constraints still face representational limitations and often underperform due to inaccurate sampling of optimal joint actions during training. To address this, we introduce an instructive

constraint to guide agents among the VD policies in the extended set  $\Pi$ :

$$\mathcal{L}_I = \min D_{\text{KL}}(\pi_\alpha^i(\cdot | \tau^i) \parallel \pi_\beta^i(\cdot | \tau^i)), \quad (8)$$

where  $\pi^i(\cdot | \tau^i)$  is a Boltzmann policy with respect to each agent’s decentralized utility function is defined as:

$$\pi_k^i(u^i | \tau^i) = \frac{\exp(Q_i(\tau^i, u^i))}{\sum_{u^i} \exp(Q_i(\tau^i, u^i))}, \forall u^i \in U^i, k \in [\alpha, \beta]. \quad (9)$$

The intuition behind this constraint is that VD methods constructing a surrogate target can learn the correct optimal joint action function, thereby identifying the optimal local behaviors. By making decentralized policies in the representationally constrained VD methods imitate this, agents are encouraged to select local optimal actions aligned with the correct joint action. This enables both VD learners to approximate the optimal joint action while adhering to the IGM principle, promoting efficient collaborative training.

#### 4.4 Overall Implementation

As shown in Figure 2, HPF consists of three parts: the joint action value function  $Q_{jt}$  estimated by the VD method with complete representation, the factored global action value function  $Q_{tot}$  and the instructive constraint  $\mathcal{L}_I$  between the utility functions of both methods. All modules are implemented by neural networks, with parameters shared across agents to promote efficient policy learning. The overall learning objective of HPF is defined as follows:

$$\mathcal{L} = \mathcal{L}_{\text{TD}}^{\text{tot}} + \mathcal{L}_{\text{TD}}^{\text{jt}} + \mathcal{L}_I. \quad (10)$$

In Equation (10),  $\mathcal{L}_{\text{TD}}^{\text{tot}}$  and  $\mathcal{L}_{\text{TD}}^{\text{jt}}$  respectively denote the update losses of the distinct heterogeneous factored global action-value function, given by:

$$\mathcal{L}_{\text{TD}}^{\text{tot}} = \sum_1^b (Q_{tot}(s, \mathbf{u}) - y'_{tot})^2, \quad (11)$$

$$\mathcal{L}_{\text{TD}}^{\text{jt}} = \sum_1^b (Q_{jt}(s, \mathbf{u}) - y'_{jt})^2, \quad (12)$$

where  $y'_{tot} = r + \gamma Q_{tot}(s', \arg\max_{\mathbf{u}'} Q_{tot}(\tau', \mathbf{u}'))$  and  $y'_{jt} = r + \gamma Q_{jt}(s', \arg\max_{\mathbf{u}'} Q_{jt}(\tau', \mathbf{u}'))$  is the respective learning target for both VD policy and  $b$  are the batches sampled by  $\Pi = [\pi_\alpha, \pi_\beta]$ .  $\mathcal{L}_I$  denotes the instructive constraint defined in Equation (8). In addition, we utilize the target networks to stabilize the training. More details of the HPF can be found in Section B of Supplementary Materials.

## 5 Experiments

In this section, we constructed two sets of HPF composite policies (WQMIX and QMIX, QPLEX and VDN, referred to as HPF-WQ and HPF-QV). Each composite policy utilizes two evaluation methods from Equation (7) for policy sampling in more difficult tasks (hence denoted as Add-HPF-WQ, Opt-HPF-WQ, Add-HPF-QV, and Opt-HPF-QV). We

$U_2 \backslash U_1$	$u_1$	$u_2$	$u_3$
$u_1$	<b>8</b>	-12	-12
$u_2$	-12	3	0
$u_3$	-12	0	3

(a) Payoff matrix

$U_2 \backslash U_1$	$u_1$	$u_2$	$u_3$
$u_1$	-8.05	-8.05	-8.06
$u_2$	-8.06	<b>1.49</b>	1.42
$u_3$	-8.06	1.44	1.37

(c) Payoff learned by QMIX

$U_2 \backslash U_1$	$u_1$	$u_2$	$u_3$
$u_1$	<b>8.03</b>	-11.98	-12.02
$u_2$	-11.98	1.53	1.54
$u_3$	-12.09	1.54	1.53

(e) Payoff learned by WQMIX

$U_2 \backslash U_1$	$u_1$	$u_2$	$u_3$
$u_1$	<b>7.99</b>	-11.99	-12.00
$u_2$	-12.02	3.00	-0.01
$u_3$	-12.04	-0.04	2.99

(g) Payoff learned by HPF-WQ

$U_2 \backslash U_1$	$u_1$	$u_3$	$u_3$
$u_1$	-6.91	-4.71	-4.62
$u_2$	-4.59	-2.41	<b>-2.33</b>
$u_3$	-4.72	-2.52	-2.46

(b) Payoff learned by VDN

$U_2 \backslash U_1$	$u_1$	$u_3$	$u_3$
$u_1$	<b>9.66</b>	-12.78	-12.76
$u_2$	-12.23	3.16	-0.75
$u_3$	-11.54	-0.70	3.12

(d) Payoff learned by QPLEX

$U_2 \backslash U_1$	$u_1$	$u_3$	$u_3$
$u_1$	-0.70	-0.24	-0.16
$u_2$	-0.26	2.01	2.52
$u_3$	-0.19	2.60	<b>3.16</b>

(f) Payoff learned by ResQ

$U_2 \backslash U_1$	$u_1$	$u_3$	$u_3$
$u_1$	<b>7.93</b>	-12.56	-12.35
$u_2$	-12.32	3.02	0.05
$u_3$	-10.67	-0.06	2.97

(h) Payoff learned by HPF-QV

Table 1: Pay-off matrix of one-step game. Boldface means optimal/greedy actions from the state-action value.

then compared these methods with state-of-the-art baselines: WQMIX, QPLEX, QMIX, VDN, and ResQ. For a fair comparison, all experimental results are illustrated with the median performance and the 25%-75% quartile over 5 random seeds. More details about hyperparameters and experiments are provided in Section C of Supplementary Materials.

### 5.1 Matrix Game

Our evaluation begins with a matrix game shown in Table 1a, where two agents need to perform the optimal joint action  $(u_1, u_1)$  to obtain the highest reward. However, sub-optimal joint actions  $(u_2, u_2)$  and  $(u_3, u_3)$  tend to make decentralized agents choose their local actions  $u_2$  or  $u_3$  when other agents select actions uniformly during training, which will lead to the noted *Relative-Overallgeneralization* problem.

Table 1 presents the comparison results of various methods in matrix game. We can observe that QMIX fails to select the optimal joint action and struggles in sub-optimal ones, while VDN even can not accurately approximate the suboptimal joint action values. Although QPLEX achieves success, it tends to overestimate the optimal joint action values. Theoretically, both WQMIX and ResQ can satisfy the IGM condition, but they still cannot correctly estimate the centralized value function distribution on this slightly more complex matrix game, and may even lead to the failure of optimal joint action selection, which may require more training. In contrast, our algorithms (HPF-WQ, HPF-QV) effectively identify and select the optimal joint action, which demonstrates superior performance than other baselines.

### 5.2 StarCraft Multi-Agent Challenge

To further showcase HPF’s scalability in more complex environments, we test it using the StarCraft Multi-Agent Chal-



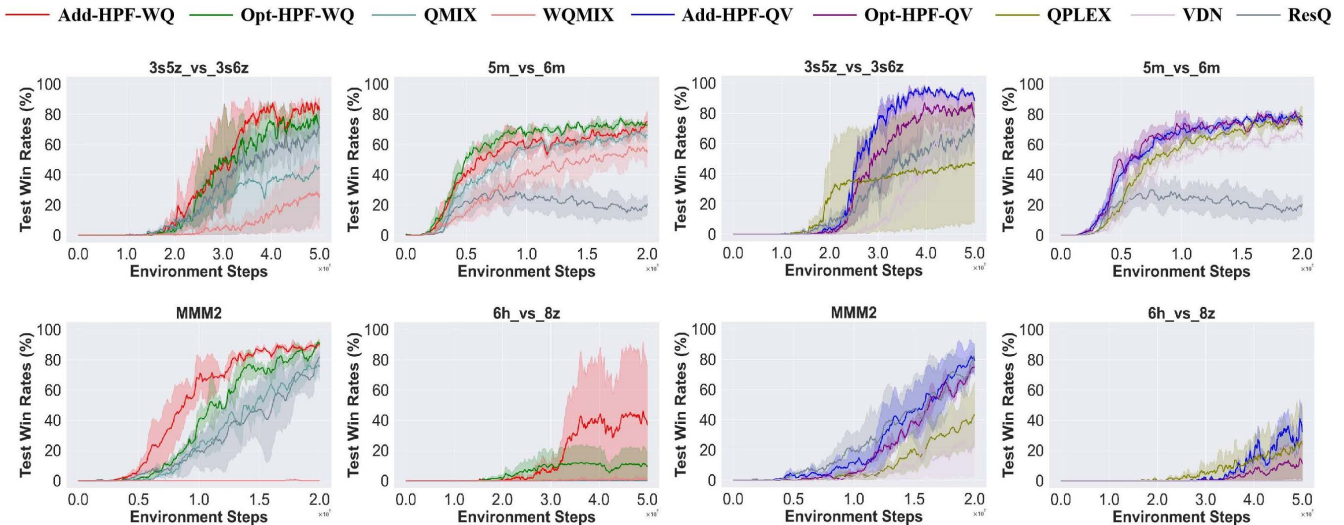


Figure 3: Comparison results on the selected scenarios in the StarCraft Multi-Agent Challenge.

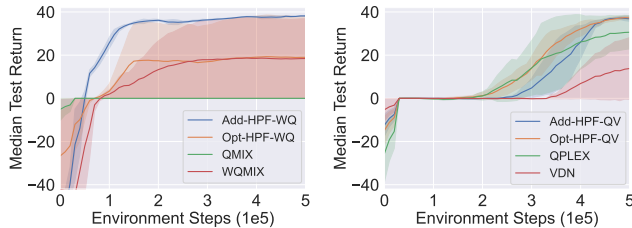


Figure 4: Comparison results in the predator and prey.

lence (SMAC) benchmark. The evaluation covers different levels of difficulty and complexity, and focuses on five maps: 3s\_vs\_5z, 5m\_vs\_6m, MMM2, 6h\_vs\_8z and 3s5z\_vs\_3s6z, which serve as our test environments.

Figure 3 shows the performance of all the algorithms on these selected scenarios. One can observe that our algorithm (Add-HPF-WQ, Opt-HPF-WQ, Add-HPF-QV and Opt-HPF-QV) significantly outperforms the counterpart baselines (WQMIX, QMIX, and QPLEX, VDN) and quickly learns a cooperative policy to achieve high testing win rates on all maps. Particularly in the super-hard challenging scenarios (6h\_vs\_8z and 3s5z\_vs\_3s6z), HPF demonstrated impressive capabilities. Though WQMIX can learn the correct action-value function, its poor performance is hard to be noted. The effectiveness of HPF stems from combining the strengths of two types of VD methods, enabling the discovery of beneficial samples during strategy sampling and interaction with the environment, thus enhancing overall performance.

### 5.3 Predator and Prey

We further evaluate the effectiveness of the HPF scheme by conducting experiments in the predator-prey environment, a partially observable task with 8 predators and 8 prey. A negative reward of -2 is given when a single predator attempts to capture prey alone, and a positive reward of +10 is received when two predators successfully cooperate to capture prey.

As shown in Figure 4, QMIX and VDN fail to learn effective policies to capture the prey. Although WQMIX benefits from the complete IGM expressiveness, its poor performance and high variance suggest difficulties in approximating the true joint action-value function. QPLEX has some competitiveness, but the effect is still inferior to our proposed scheme, which further demonstrates the effectiveness of HPF.

### 5.4 Ablation Studies

To test the effectiveness of the sampling method for VD policies in HPF, we replaced the sampling approach in Equations (5)-(6) with random sampling to validate its effectiveness. The experimental results, as presented in Section D of Supplementary Materials, demonstrate that the composite policy in HPF with sampling based on the estimated value function achieves superior performance. The sampling method in HPF enables the model to identify more suitable sampling policies to interact, thus improving the overall performance of the model.

## 6 Conclusion

This paper introduces the Heterogeneous Policy Fusion scheme for cooperative multi-agent reinforcement learning, which balances representation ability and training efficiency in value decomposition algorithms by combining policies with both restricted and unrestricted representation capacities. HPF designs a composite policy by extending two distinct types of value decomposition policies into a policy set, and integrating the interaction policy based on their value function estimates. This extension enables HPF to attain both high training efficiency and complete representation capacity. Experiments show that HPF not only learns optimal joint actions effectively but also outperforms most baseline VD methods without needing specialized design, proving its effectiveness and ease of use.

## References

- [Ding *et al.*, 2023a] Shifei Ding, Wei Du, Ling Ding, Lili Guo, Jian Zhang, and Bo An. Multi-agent dueling q-learning with mean field and value decomposition. *Pattern Recognition*, 139:109436, 2023.
- [Ding *et al.*, 2023b] Shifei Ding, Wei Du, Ling Ding, Jian Zhang, Lili Guo, and Bo An. Multiagent reinforcement learning with graphical mutual information maximization. *IEEE Transactions on Neural Networks and Learning Systems*, 2023.
- [Jaques *et al.*, 2019] Natasha Jaques, Angeliki Lazaridou, Edward Hughes, Caglar Gulcehre, Pedro Ortega, DJ Strouse, Joel Z Leibo, and Nando De Freitas. Social influence as intrinsic motivation for multi-agent deep reinforcement learning. In *International conference on machine learning*, pages 3040–3049. PMLR, 2019.
- [Khan *et al.*, 2022] Muhammad Junaid Khan, Syed Hammad Ahmed, and Gita Sukthankar. Transformer-based value function decomposition for cooperative multi-agent reinforcement learning in starcraft. In *Proceedings of the AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment*, volume 18, pages 113–119, 2022.
- [Liu *et al.*, 2023] Shunyu Liu, Yihe Zhou, Jie Song, Tongya Zheng, Kaixuan Chen, Tongtian Zhu, Zunlei Feng, and Mingli Song. Contrastive identity-aware learning for multi-agent value decomposition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 11595–11603, 2023.
- [Mei *et al.*, 2023] Yongsheng Mei, Hanhan Zhou, Tian Lan, Guru Venkataramani, and Peng Wei. Mac-po: Multi-agent experience replay via collective priority optimization. In *Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent Systems*, pages 466–475, 2023.
- [Mnih *et al.*, 2015] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Belle-mare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, and Georg Ostrovski. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- [Naderializadeh *et al.*, 2020] Navid Naderializadeh, Fan H Hung, Sean Soleyman, and Deepak Khosla. Graph convolutional value decomposition in multi-agent reinforcement learning. *arXiv preprint arXiv:2010.04740*, 2020.
- [Oliehoek and Amato, 2016] Frans A. Oliehoek and Christopher Amato. *A Concise Introduction to Decentralized POMDPs*. Springer Briefs in Intelligent Systems. Springer, 2016.
- [Perrusquía *et al.*, 2021] Adolfo Perrusquía, Wen Yu, and Xiaou Li. Multi-agent reinforcement learning for redundant robot control in task-space. *International Journal of Machine Learning and Cybernetics*, 12:231–241, 2021.
- [Pina *et al.*, 2022] Rafael Pina, Varuna De Silva, Joosep Hook, and Ahmet Kondo. Residual q-networks for value function factorizing in multiagent reinforcement learning. *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- [Rashid *et al.*, 2018] Tabish Rashid, Mikayel Samvelyan, Christian Schröder de Witt, Gregory Farquhar, Jakob N. Foerster, and Shimon Whiteson. QMIX: Monotonic Value Function Factorisation for Deep Multi-Agent Reinforcement Learning. In *Proceedings of the 35th International Conference on Machine Learning*, pages 4292–4301. PMLR, 2018.
- [Rashid *et al.*, 2020] Tabish Rashid, Gregory Farquhar, Bei Peng, and Shimon Whiteson. Weighted qmix: Expanding monotonic value function factorisation for deep multi-agent reinforcement learning. *Advances in neural information processing systems*, 33:10199–10210, 2020.
- [Samvelyan *et al.*, 2019] Mikayel Samvelyan, Tabish Rashid, Christian Schröder de Witt, Gregory Farquhar, Nantas Nardelli, Tim G. J. Rudner, Chia-Man Hung, Philip H. S. Torr, Jakob N. Foerster, and Shimon Whiteson. The StarCraft Multi-Agent Challenge. In *Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent Systems, AAMAS*, pages 2186–2188, 2019.
- [Shen *et al.*, 2022] Siqi Shen, Mengwei Qiu, Jun Liu, Weiquan Liu, Yongquan Fu, Xinwang Liu, and Cheng Wang. Resq: A residual q function-based approach for multi-agent reinforcement learning value factorization. *Advances in Neural Information Processing Systems*, 35:5471–5483, 2022.
- [Son *et al.*, 2019] Kyunghwan Son, Daewoo Kim, Wan Ju Kang, David Earl Hostallero, and Yung Yi. Qtran: Learning to factorize with transformation for cooperative multi-agent reinforcement learning. In *International conference on machine learning*, pages 5887–5896. PMLR, 2019.
- [Son *et al.*, 2020] Kyunghwan Son, Sungsoo Ahn, Roben Delos Reyes, Jinwoo Shin, and Yung Yi. Qtran++: improved value transformation for cooperative multi-agent reinforcement learning. *arXiv preprint arXiv:2006.12010*, 2020.
- [Sunehag *et al.*, 2017] Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi, Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z Leibo, Karl Tuyls, et al. Value-decomposition networks for cooperative multi-agent learning. *arXiv preprint arXiv:1706.05296*, 2017.
- [Sutton and Barto, 2018] Richard S. Sutton and Andrew G. Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [Tan, 1993] Ming Tan. Multi-agent reinforcement learning: Independent vs. cooperative agents. In *Proceedings of the 10th international conference on machine learning*, pages 330–337. PMLR, 1993.
- [Wang *et al.*, 2021] Jianhao Wang, Zhizhou Ren, Terry Liu, Yang Yu, and Chongjie Zhang. {QPLeX}: Duplex dueling multi-agent q-learning. In *International Conference on Learning Representations*, 2021.



- [Wei *et al.*, 2022] Qinglai Wei, Yugu Li, Jie Zhang, and Fei-Yue Wang. Vgn: Value decomposition with graph attention networks for multiagent reinforcement learning. *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- [Xu *et al.*, 2023] Zhiwei Xu, Yunpeng Bai, Bin Zhang, Dapeng Li, and Guoliang Fan. Haven: Hierarchical cooperative multi-agent reinforcement learning with dual coordination mechanism. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 11735–11743, 2023.
- [Yang *et al.*, 2020] Yaodong Yang, Jianye Hao, Ben Liao, Kun Shao, Guangyong Chen, Wulong Liu, and Hongyao Tang. Qatten: A general framework for cooperative multiagent reinforcement learning. *arXiv preprint arXiv:2002.03939*, 2020.
- [Yang *et al.*, 2023] Shantian Yang, Bo Yang, Zheng Zeng, and Zhongfeng Kang. Causal inference multi-agent reinforcement learning for traffic signal control. *Information Fusion*, 94:243–256, 2023.
- [Zhang *et al.*, 2023] Ruilong Zhang, Qun Zong, Xiuyun Zhang, Liqian Dou, and Bailing Tian. Game of drones: Multi-uav pursuit-evasion game with online motion planning by deep reinforcement learning. *IEEE Transactions on Neural Networks and Learning Systems*, 34(10):7900–7909, 2023.
- [Zohar *et al.*, 2022] Roy Zohar, Shie Mannor, and Guy Tenenbholz. Locality matters: A scalable value decomposition approach for cooperative multi-agent reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 9278–9285, 2022.

## 7 Appendix: Brief Introduction

This supplementary material provides the proof procedure, additional environment settings, and ablation results that were omitted from the main text. In Section A, we provide the proof of Proposition 1. Section B describes the specific implementation process of HPF, which takes WQMIX and QMIX as examples of two candidate VD policies of different types. We describe the experimental scenarios and hyperparameters setup in Section C in detail. The additional analysis and omitted ablation results are presented in Section D.

### A Proof of Proposition 1

In this section, we show the proof of Proposition 1.

**Proposition 2.** *If  $\Pi = [\pi_\alpha, \pi_\beta]$  is an extended policy set formed by existing heterogeneous VD policies, then  $\Pi$  still satisfies the IGM criterion.*

*Proof.* The Theorem 1 and Section 4 of ResQ [Shen *et al.*, 2022] show that for any hard-to-factorize state-action value function  $Q_{jt}$ . It can generally always be regarded as satisfying the following formula:

$$Q_{jt}(\tau, \mathbf{u}) = w_{tot}(\tau, \mathbf{u})Q_{tot}(\tau, \mathbf{u}) + w_r(\tau, \mathbf{u})Q_r(\tau, \mathbf{u}) \quad (13)$$

where masks  $w_{tot}(\tau, \mathbf{u}), w_r(\tau, \mathbf{u}) \in \{0, 1\}$ . The main  $Q_{tot}$  shares the same greedy optimal policy as  $Q_{jt}$ , and  $Q_{tot}(\tau, \mathbf{u})$  is factorized into per-agent utilities  $Q_i$ , and each agent selects its own optimal action greedily with respect to  $Q_i$ .

More generally, the value decomposition policies that construct surrogate targets can all be regarded as special cases of Formula (13) and can be represented in a generalized form as:

$$Q_{jt}^{Type1}(\tau, \mathbf{u}) = Q_{tot}^{(1)}(\tau, \mathbf{u}) + w_r(\tau, \mathbf{u})Q_r(\tau, \mathbf{u}) \quad (14)$$

where  $Q_r(\tau, \mathbf{u}) \leq 0$ ,  $Q_{tot}^{(1)}(\tau, \mathbf{u})$  and  $[Q_i(\tau_i, u_i)]_{i=1}^n$  satisfy the monotonicity condition  $\partial Q_{tot}^{(1)}(\tau, \mathbf{u}) / \partial Q_i(\tau_i, u_i) \geq 0, \forall i \in N$  and

$$w_r(\tau, \mathbf{u}) = \begin{cases} 0 & \mathbf{u} = \bar{\mathbf{u}} \\ 1 & \mathbf{u} \neq \bar{\mathbf{u}} \end{cases} \quad (15)$$

While the value decomposition policies by network constraints still satisfy the monotonicity and are depicted as:

$$Q_{jt}^{Type2}(\tau, \mathbf{u}) = Q_{tot}^{(2)}(\tau, \mathbf{u}) = \sum_{i=1}^n |w_i| Q_i(\tau_i, u_i) \quad (16)$$

For the extended composite policy set  $\Pi = [\pi_\alpha, \pi_\beta]$  in the HPF framework, the VD policies it contains are sampled for interaction and update, so its centralized value function  $Q_{jt}$  is expressed as:

$$\begin{aligned} Q_{jt}(\tau, \mathbf{u}) &= \lambda Q_{jt}^{Type1}(\tau, \mathbf{u}) + (1 - \lambda) Q_{jt}^{Type2}(\tau, \mathbf{u}) \\ &= \lambda(Q_{tot}^{(1)}(\tau, \mathbf{u}) + w_r(\tau, \mathbf{u})Q_r(\tau, \mathbf{u})) \\ &\quad + (1 - \lambda)Q_{tot}^{(2)}(\tau, \mathbf{u}) \end{aligned} \quad (17)$$

where  $Q_r(\boldsymbol{\tau}, \mathbf{u}) \leq 0$  and  $\lambda \in [0, 1]$  is the sample probability of the candidate VD policy in  $\Pi$ . We show that  $\arg \max_{\mathbf{u}} Q_{jt}(\boldsymbol{\tau}, \mathbf{u}) = \bar{\mathbf{u}}$ , where  $\bar{u}_i = \arg \max_{u_i} Q_i(\tau_i, u_i)$  and  $\bar{\mathbf{u}} = [\bar{u}_i]_{i=1}^n$ .

Then the inequality holds:

$$Q_{jt}(\boldsymbol{\tau}, \bar{\mathbf{u}}) = \lambda Q_{tot}^{(1)}(\boldsymbol{\tau}, \bar{\mathbf{u}}) + (1 - \lambda) Q_{tot}^{(2)}(\boldsymbol{\tau}, \bar{\mathbf{u}}) \quad (18a)$$

$$\geq \lambda Q_{tot}^{(1)}(\boldsymbol{\tau}, \mathbf{u}) + (1 - \lambda) Q_{tot}^{(2)}(\boldsymbol{\tau}, \mathbf{u}) \quad (18b)$$

$$\geq \lambda (Q_{tot}^{(1)}(\boldsymbol{\tau}, \mathbf{u}) + w_r(\boldsymbol{\tau}, \mathbf{u}) Q_r(\boldsymbol{\tau}, \mathbf{u})) + (1 - \lambda) Q_{tot}^{(2)}(\boldsymbol{\tau}, \mathbf{u}) \quad (18c)$$

$$= Q_{jt}(\boldsymbol{\tau}, \mathbf{u}) \quad (18d)$$

(18c) comes from  $w_r(\boldsymbol{\tau}, \mathbf{u}) = 1, \forall \mathbf{u} \neq \bar{\mathbf{u}}$  and  $Q_r(\boldsymbol{\tau}, \mathbf{u}) \leq 0$ . This inequality (18) means that  $\bar{\mathbf{u}} = [\bar{u}_i]_{i=1}^n$  maximizes  $Q_{jt}(\boldsymbol{\tau}, \mathbf{u})$  with both  $Q_{tot}^{(2)}(\boldsymbol{\tau}, \mathbf{u})$  and  $Q_{tot}^{(1)}(\boldsymbol{\tau}, \mathbf{u})$  satisfy the monotonicity condition. Thus  $[Q_i(\tau_i, u_i)]_{i=1}^n$  satisfies the IGM for  $Q_{jt}(\boldsymbol{\tau}, \mathbf{u})$ .  $\square$

## B Detailed Implementation of HPF

In this section, we show the explicit implementation of HPF, which is also described in Algorithm 1. Here we take Add-HPF-WQ as an example to illustrate the HPF algorithm process, with other HPF variants following a similar training procedure.

In Add-HPF-WQ, WQMIX[Rashid *et al.*, 2020] is chosen as  $\pi_\alpha$  in Eq.(4) of the main text. While the  $\pi_\beta$  with limited representation is QMIX[Rashid *et al.*, 2018]. Since QMIX assumes the task can be directly decomposed, which also means  $Q_{jt} = Q_{tot}$  is this setting.

First, we initialize the network parameters of the two types of value decomposition policy  $\theta_\alpha, \theta_\beta$  as well as their target network parameters  $\theta_\alpha^-, \theta_\beta^-$ , experience replay buffer  $\mathcal{D}$ , learning rate  $\zeta$ , the size of batch sampling *batch\_size* and target network *update\_interval*  $m$ . We evaluate the utility functions of both candidate VD policies,  $[Q_i(\tau^i, u^i; \theta_\alpha)]_{i=1}^n$  and  $[Q_i(\tau^i, u^i; \theta_\beta)]_{i=1}^n$ , to construct the sampling policy for composite  $\Pi$  via Eq.(5)-Eq.(7) in the main text. Then the interactive trajectory will be stored in the replay buffer.

During the training, HPF extracts a small batch of interaction trajectories from  $\mathcal{D}$  and estimates the utility function  $[Q_i^{\pi_\alpha}]_{i=0}^n, [Q_i^{\pi_\beta}]_{i=0}^n$  and the centralized value function  $Q_{jt}(\theta_\alpha), Q_{jt}(\theta_\beta)$  for each time step. The *Mlp*( $\cdot$ ) function of Algorithm 1 represents a regular feedforward neural network that does not manipulate the weight parameters generated by the hypernetwork in any way. While *Mix*( $\cdot$ ) function represents the mixing network in QMIX, which involves taking the absolute values of the weight generated by the hypernetwork, ensuring all weights remain positive to meet IGM principles. Then we can update  $Q_{jt}(\theta_\alpha), Q_{jt}(\theta_\beta)$  and all the utility functions like DQN[Mnih *et al.*, 2015] updates through the back-propagation. In the following part, we omit the parameters in the functions, such as  $Q_{jt}^{\pi_\alpha} = Q_{jt}^{\pi_\alpha}(\theta_\alpha), Q_{jt}^{\pi_\alpha(-)} = Q_{jt}^{\pi_\alpha}(\theta_\alpha^-)$ , and  $Q_a^{\pi_\alpha} = Q_a(\tau_t^a, u_t^a), a \in \{1, \dots, n\}$  for brevity.

---

### Algorithm 1 Add-HPF-WQ for MARL

---

**Require:**  $\theta_\alpha, \theta_\beta, \theta_\alpha^-, \theta_\beta^-$

**Initialize**  $\theta_\alpha^- \leftarrow \theta_\alpha, \theta_\beta^- \leftarrow \theta_\beta, \zeta, \mathcal{D}, b, \text{step}=0$ , update *interval*= $m$

```

1: while step < step_max do
2:    $t = 0, s_0 = \text{initial state}$ ;
3:   while  $s_t \neq \text{terminal}$  and  $t < \text{episode limit}$  do
4:     for each agent  $a$  do
5:        $\tau_t^a = \tau_{t-1}^a \cup \{(o_t, u_{t-1})\}$ ;
6:       Select action  $u_t^a$  by  $\epsilon$ -greedy scheme;
7:       Get the utilities  $Q_a^{\pi_\alpha}, Q_a^{\pi_\beta}, a \in \{1, \dots, n\}$ ;
8:       Obtain  $Q^{\pi_\alpha}$  and  $Q^{\pi_\beta}$  via Eq.(7) in the main text
9:       Calculate  $P_w$  and sampling  $\mathbf{u}_t$  via Eq.(5)-Eq.(6)
           in the main text
10:    end for
11:    Get reward  $r_t$  and next state  $s_{t+1}$ ;
12:     $\mathcal{D} = \mathcal{D} \cup \{(s_t, \mathbf{u}_t, r_t, s_{t+1})\}; t = t + 1$ ;
13:    step = step + 1;
14:  end while
15:  if  $|\mathcal{D}| > \text{batch-size}$  then
16:     $b \leftarrow \text{select random batch of episodes from } \mathcal{D}$ ;
17:    for  $t$  in each trajectory in batch  $b$  do
18:       $Q_{jt}^{\pi_\alpha} = \text{Mlp}(Q_1^{\pi_\alpha}, \dots, Q_n^{\pi_\alpha}, s_t)$ ;
19:       $Q_{tot}^{\pi_\beta} = \text{Mix}(Q_1^{\pi_\beta}, \dots, Q_n^{\pi_\beta}, s_t)$ ;
20:      Estimate targets  $Q_{jt}^{\pi_\alpha(-)}$  and  $Q_{tot}^{\pi_\beta(-)}$ ;
21:    end for
22:    Update the joint action value function:
23:     $\Delta Q_{jt}^{\pi_\alpha} = Q_{jt}^{\pi_\alpha(-)} - Q_{jt}^{\pi_\alpha}$ ,
24:     $\Delta \theta_\alpha = \nabla_{\theta_\alpha} (\Delta Q_{jt}^{\pi_\alpha})^2, \theta_\alpha = \theta_\alpha - \zeta \Delta \theta_\alpha$ ;
25:     $\Delta Q_{tot} = Q_{tot}^{\pi_\beta(-)} - Q_{tot}^{\pi_\beta}$ ,
26:     $\Delta \theta_\beta = \nabla_{\theta_\beta} (\Delta Q_{tot}^{\pi_\beta})^2, \theta_\beta = \theta_\beta - \zeta \Delta \theta_\beta$ ;
27:  end if
28:  if step %  $m = 0$  then
29:    Update the target network  $\theta_\alpha^- \leftarrow \theta_\alpha, \theta_\beta^- \leftarrow \theta_\beta$ ;
30:  end if
31: end while

```

---

## C Experiment and Hyperparameters Setting

In this section, we thoroughly describe the experimental setups and parameter configurations for the algorithms employed in this paper. The following subsections elaborate on the specific objectives of each experimental scenario and the neural network parameters of the compared algorithms.

### C.1 Description of the experimental scenarios

**Matrix Game.** The matrix game is a simple mathematical model used to study the interaction and decision-making scenarios of two agents in a shared environment. In this setting, each agent has three possible actions, and their choices result in simultaneous outcomes according to a shared payoff function. This symmetric matrix game model exhibits an optimal joint action  $(u_1, u_1)$ , where agents possess an observation with all values set to 1, indicating that agents solely consider interactive actions to impact specific cooperative rewards. Agents can only choose one action per round of inter-

action. This model allows for a rapid evaluation of whether the tested algorithms satisfy the IGM condition and whether the centralized value function under the joint actions of agents can accurately estimate the value functions.

**Predator-Prey.** The *relative-overgeneralization* is formulated by a predator-prey task that is modeled in a partially observable grid-world: 8 agents try to coordinate to hunt 8 prey in a  $10 \times 10$  grid. Each agent possesses a  $5 \times 5$  sub-grid sight range. The prey will be caught if two adjacent predators are executing the *catch* action. The predators are rewarded  $r = 10$  when they capture prey successfully but get punished reward  $p = -2$  when only a single agent attempts to capture any prey. Collaboration among the predators is essential to capture all prey successfully and prevent incurring penalties.

**StarCraft Multi-Agent Challenge.** StarCraft II, being a real-time strategy game, presents an excellent opportunity to address cooperative challenges in the multi-agent domain. SMAC [Samvelyan *et al.*, 2019] utilizes Blizzard’s StarCraft II Machine Learning API and DeepMind’s PySC2 as an interface to engage with the game environment. Our focus lies on the micromanagement challenges within SMAC, where each unit is under the control of an independent agent, relying solely on limited local observation. These unit groups must be trained to engage an opposing army under the centralized control of the game’s built-in scripted AI. The objective of cooperative allies is to maximize damage dealt to enemy units while minimizing received damage. Symmetrically placed opposing groups with varied initial unit placements characterize the scenario across episodes.

## C.2 Hyperparameters Setting

For the sake of consistency, we integrate all algorithms in PyMARL and apply the same network architecture and hyperparameters to the contrasting methods. In alignment with PyMARL, the agent architecture comprises a DRQN with a recurrent layer, featuring a 64-dimensional hidden state. Throughout the entire training process, each agent employs an independent  $\epsilon$ -greedy exploration action selection based on its own  $Q_i$ . The  $\epsilon$  value is annealed from 1.0 to 0.05 over  $50k$  time steps and remains constant thereafter ( $100k$  for super-hard scenarios 6h\_vs\_8z and 3s5z\_vs\_3s6z). The discount factor,  $\gamma = 0.99$ , remains constant across all experiments. RMSprop with  $\alpha = 0.99$  is used without incorporating weight decay or momentum. The learning rate for training is consistently set to  $5 \times 10^{-4}$  in all testing scenarios in Predator-Prey and baselines in SMAC. We apply layer-normalization to the observations, and to ensure stability when updating the composite policy  $\Pi = [\pi_\alpha, \pi_\beta]$  with sampled trajectories, the Adam optimizer was utilized in the HPF scheme with the default setting, with the temperature  $\eta = 1$  in candidate policies sampling in HPF.<sup>1</sup>

As for The Matrix Game, all agents adhere to full exploration, i.e.,  $\epsilon = 1$ , and remain unchanged during training. The WQMIX baseline and the WQMIX policy in HPF employ a sub-optimal joint action weight of  $w = 0.1$ , while in the SMAC environment, it opts for  $w = 0.75$  from the range

$w = (0.5, 0.75)$ . While the QPLEX baseline and the QPLEX policy in HPF enjoy the same setting as in its open-source code for fair comparison.

Training halts every 10000-time steps for 16 evaluation episodes. The most recent 5000 training episodes are maintained in the replay buffer, and batches of 32 episodes are uniformly sampled for training. All agent networks share the parameters of networks, and we update the target networks every 200 episodes. To ensure fairness in comparing the performance, all other parameters of the baselines are consistent with their corresponding settings in the official code bases.

## D Additional Experimental Results

In this section, we supplement the ablation experiments omitted in the main text. Specifically, we replace the sampling method of the composite policy in the HPF scheme with random sampling to verify the importance of utilizing value es-

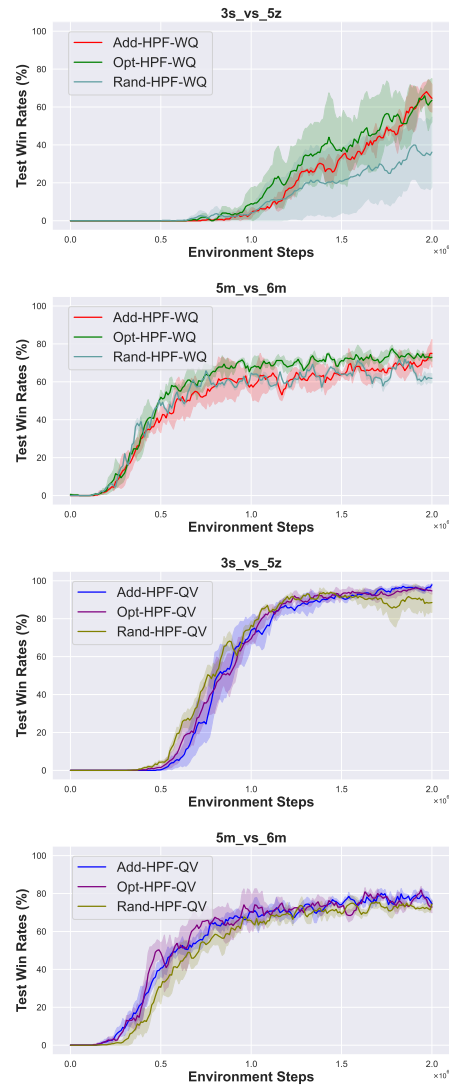


Figure 5: Ablation studies of the random candidate VD policy sampling in HPF.

<sup>1</sup>We use the SC2.4.10 version of SMAC through all the testing scenarios.

estimates in VD methods as sampling weights in HPF. Additionally, we record the sampling frequency of heterogeneous VD methods within HPF across different scenarios, attempting to explain which type of VD method played a greater role in model updates.

### D.1 Omitted Ablation Studies

This subsection provides the ablation experimental results of a random sampling of the heterogeneous VD method in the HPF scheme. As illustrated in Figure 5, when HPF adopts random sampling to create a composite policy, its performance in the 3s\_vs\_5z and 5m\_vs\_6m scenarios is inferior to the approach proposed in this paper, which utilizes the sampling from a Boltzmann policy based on value function

estimates of the heterogeneous VD method, as depicted by Eq.(5)-Eq.(6) in the main text. Since a greater estimate implies that the currently sampled VD method has a higher cumulative expected return, which makes the data obtained by the interaction of this method with the environment have a higher potential value for the model update. This indicates that the composite policy of HPF is capable of effectively integrating the merits of the heterogeneous VD method and thereby acquiring better interaction experiences, consequently enhancing the training efficiency of the model.

### D.2 Additional Analysis

The main purpose of this subsection is to analyze which type of VD policy is more likely to be chosen in HPF, providing insights into the impact of interaction data generated by different VD methods on the model. We recorded the probability of sampling a certain type of VD policy during model training in Add-HPF-WQ, as illustrated in Figure 6. Here, *wqmix\_policy* represents the probability of choosing the WQMIX policy, while *qmix\_policy* represents the probability of choosing the QMIX policy. The selection of VD policies by the HPF scheme varies across different cooperative tasks. However, a general observation is that in tasks where the WQMIX policy consistently learns effective collaborative behaviors, such as 5m\_vs\_6m, HPF tends to choose and maintain the WQMIX policy. While in the tasks where obtaining effective information from WQMIX interaction experiences proves challenging, HPF gradually leans towards selecting the QMIX policy for interaction, thereby enhancing the training effectiveness of the model. For instance, in the case of 3s\_vs\_5z, the probability of choosing the WQMIX policy increases at the beginning. Still, as it fails to learn suitable cooperative behaviors, HPF gradually shifts towards selecting the QMIX policy for interaction. In more complex scenarios such as MMM2, HPF progressively increases the probability of selecting the QMIX policy for interaction. This suggests that the HPF’s approach of extending VD methods can effectively perceive the task complexity of the environment and adaptively adjust the sampling probability of interaction, enabling the model to learn superior collaborative policies.

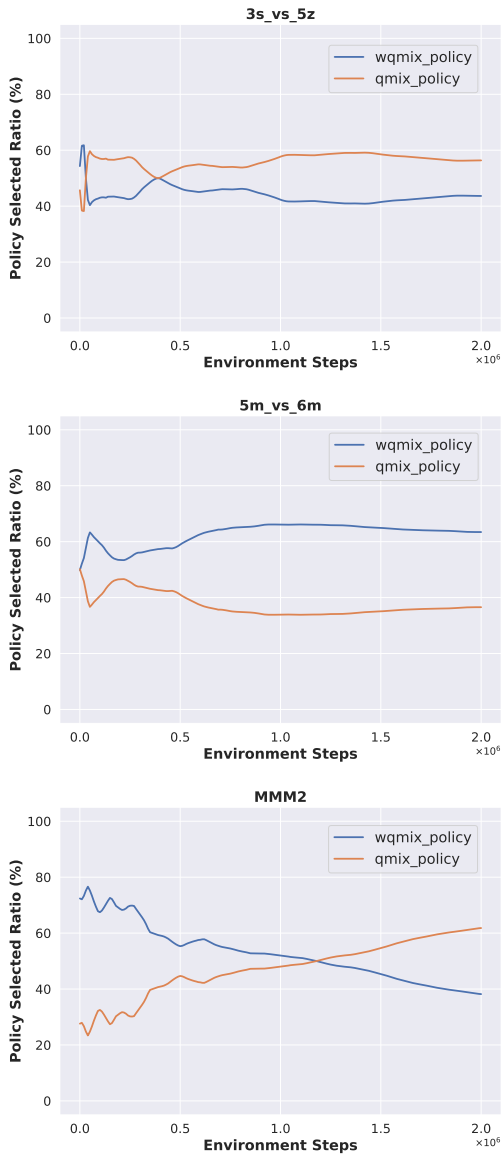


Figure 6: The selection probabilities of different VD policies in HPF along with training process.