# A Survey of Multi Agent Reinforcement Learning Federated Learning and Cooperative and Noncooperative Decentralized Regimes

Kemboi Cheruiyot, Nickson Kiprotich, Vyacheslav Kungurtsev,
Kennedy Mugo, Vivian Mwirigi, Marvin Ngesa

July 2025

### Abstract

The increasing interest in research and innovation towards the development of autonomous agents presents a number of complex yet important scenarios of multiple AI Agents interacting with each other in an environment. The particular setting can be understood as exhibiting three possibly topologies of interaction - centrally coordinated cooperation, ad-hoc interaction and cooperation, and settings with noncooperative incentive structures. This article presents a comprehensive survey of all three domains, defined under the formalism of Federal Reinforcement Learning (RL), Decentralized RL, and Noncooperative RL, respectively. Highlighting the structural similarities and distinctions, we review the state of the art in these subjects, primarily explored and developed only recently in the literature. We include the formulations as well as known theoretical guarantees and highlights and limitations of numerical performance. [1]

## 1 Introduction

A process of significant contemporary interest is multiple autonomous AI systems in the same environment as described by Multi-Agent Reinforcement Learning (MARL). There are three distinct circumstances that define particular regimes of interaction: 1) a collection of agents with a means of centralizing communication that are learning about the environment together, defined as **Federated RL**, 2) a swarm of agents with a network structure of gossip peer-to-peer communication cooperating towards common and individual goals, defined as **Decentralized MARL** and 3) interacting AI agents within the context of a noncooperative game, which we denote as **Noncooperative MARL**.

This Survey presents the state of the art in this field that has largely been developing over, approximately, the last five years. As General Purpose AI become increasingly more capable, Agents engage with an increasing number of administrative computer tasks, and robots and drones become operational in a wider breadth of physical circumstances, understanding the training for cooperation and strategic decision making becomes increasingly critical for ensuring their operation is well-aligned to be safe, effective and robust.

In the broader context of AI Alignment, the additional moving parts of multiple agents presents a greater challenge to establishing forward and backward alignment of user and stakeholder intentions. In addition, potential conflict among two AI Systems corresponding to distinct users can result in unpredictable problems and consequences. Understanding these systems and developing reliable training and communication algorithms will be an active area of research to accommodate the developing multi-agent AI system industry.

We proceed, in the next Section, to provide a brief overview of the main concepts of Reinforcement Learning, We then proceed with our survey. For each class of architecture, we present their form of deployment and structure of communication, the Markov Decision Process (MDP) nomenclature definition of the setting, present formal algorithms representative of the literature, discuss theoretical guarantees, and describe open problems in active or future pursuit by the respective research community.

- **Federated Reinforcement Learning (FRL)**: This constitutes distributed computing client teamwork in solving RL problems, where the clients can reside on distinct autonomous systems. Typically, a central server relays communication between clients and coordinates their training, although all-to-all communication as centralization is also a possible architecture. This involves multiple agents learning locally, keeping data decentralized to preserve privacy of the client. Parameters are exchanged with the eerver after sufficient epochs of local training and agglomerated with some averaging operator by the server, which then relays the centralized parameters back across the clients for the next round of training.

- **Decentralized Multi-Agent Reinforcement Learning (DMARL)**: Here, agents cooperate without a central server, communicating directly peer-to-peer. The structure of available communication channels is defined by a graph, which can be directed or undirected. This is also known as a Gossip network. They cooperate by communicating parameters in a manner so as to encourage diffusion across the network, while aiming to maximize common and/or individual rewards.

- **Noncooperative Multi-Agent Reinforcement Learning (NMARL)**: Alternatively called *Nash MARL*, this defines a set of coupled RL problems wherein the reward function for each agent is distinct and depends on the actions of all of the agents. Typically, there are competitive trade-offs in the reward structure of the states, with few relevant states that are universally high, and few that are universally low. Understanding

2

and computing various forms of Nash equilibria are described from the literature.

# 2  Introduction to Reinforcement Learning

Reinforcement Learning (RL) is a branch of machine learning that focuses on sequential decision-making, where an agent interacts with an environment to learn an optimal policy through trial and error. Unlike supervised learning, which relies on labeled data, RL involves learning from delayed feedback, where the agent receives rewards or penalties based on its actions, see Sutton and Barto (2018) for a classic text. A key challenge in RL is balancing exploration—where the agent tries new actions to discover better strategies—with exploitation, where it leverages known information to maximize immediate rewards. Striking the right balance between these two aspects is crucial for long-term success in uncertain environments .

Over the decades, RL has evolved from early studies in behavioral psychology and trial-and-error learning to form the backbone of modern algorithms that have achieved remarkable success across diverse applications. High-profile demonstrations, such as the triumph of AlphaGo in defeating world champions and significant advancements in robotics and autonomous driving, highlight the potential of RL to transform complex, real-world decision-making problems Naeem et al. (2020).

At its core, RL is built upon key concepts such as states, actions, rewards, and policies, underpinned by the mathematical framework of Markov Decision Processes (MDPs). This theoretical foundation not only enables the development of robust algorithms but also highlights inherent challenges, such as sample efficiency, scalability, and the delicate balance between exploration and exploitation.

Recent advances, particularly in deep reinforcement learning, have further expanded the capabilities of RL by integrating deep neural networks with traditional RL methods. This fusion has opened up new avenues for solving problems in high-dimensional spaces and complex environments, while ongoing research continues to address issues like safety, stability, and multi-agent interactions Gavi (2025).

This section provides a comprehensive overview of Reinforcement Learning, the fundamentals, key algorithms and methods.

## 2.1  Markov Decision Processes

The formal foundation of RL is established through the Markov Decision Process (MDP), a mathematical model for decision-making under uncertainty. An MDP describes an environment where an agent occupies a state, selects an action, transitions to a new state, and receives a corresponding reward. The agent's objective is to determine an optimal policy that maximizes the expected

cumulative discounted reward, ensuring that both immediate and future consequences of actions are considered. The framework of MDPs is central to RL, as it provides a structured way of modeling decision-making problems. An MDP is formally defined as a tuple $(S, A, P, R, \gamma)$, where:

- $S$ is the set of states,

- $A$ is the set of actions available to the agent,

- $P(s'|s, a)$ is the transition probability from state $s$ to state $s'$ given action $a$,

- $R(s, a)$ is the reward function that provides feedback to the agent,

- $\gamma \in [0, 1]$ is the discount factor, which determines the importance of future rewards.

The dynamics of MDPs are captured by the transition probabilities $P(s'|s, a)$, which define the likelihood of transitioning from one state to another given a particular action. The reward function $R(s, a)$ specifies the immediate benefit of taking an action in a given state, while the discount factor $\gamma$ balances the trade-off between immediate and future rewards.

A key problem in RL is to find an *optimal policy* $\pi^*$ that prescribes the best action for each state to maximize the expected cumulative reward. The value function $V^\pi(s)$ measures the expected return starting from state $s$ under policy $\pi$ and is given by:

$$V^\pi(s) = \sum_{a \in A} \pi(a|s) \left[ R(s, a) + \gamma \sum_{s'} P(s'|s, a) V^\pi(s') \right].$$

This formulation clearly shows the expectation over both the action selection and the state transition.

## 2.2 Dynamic Programming and Bellman Equations

Dynamic Programming (DP) algorithms solve MDPs by decomposing a complex problem into simpler sub-problems. The Bellman equations are at the heart of DP methods, providing a recursive relationship that expresses the value of a state (or action) in terms of the values of subsequent states.

**Value Iteration**  Value Iteration computes the optimal value function $V^*(s)$ by iteratively applying the Bellman optimality equation:

$$V^*(s) = \max_{a \in A} \left[ R(s, a) + \gamma \sum_{s'} P(s'|s, a) V^*(s') \right].$$

This equation indicates that the value of a state is the maximum expected return achievable by considering all possible actions and future states.

**Policy Iteration**  Policy Iteration involves two main steps: policy evaluation and policy improvement.

**Policy Evaluation:**

$$V^\pi(s) = \sum_{a \in A} \pi(a|s) \left[ R(s,a) + \gamma \sum_{s'} P(s'|s,a) V^\pi(s') \right].$$

**Policy Improvement:**

$$\pi'(s) = \arg \max_{a \in A} \left[ R(s,a) + \gamma \sum_{s'} P(s'|s,a) V^\pi(s') \right].$$

Although DP methods guarantee convergence, they require complete knowledge of the environment (i.e., the transition probabilities $P(s'|s,a)$), which is often impractical in real-world applications.

## 2.3   Model-Based and Model-Free Approaches

Reinforcement learning methods are broadly categorized based on whether they utilize an explicit model of the environment. Model-based methods incorporate or learn a representation of the environment's dynamics, enabling planning via simulated trajectories. This approach often results in superior sample efficiency since the agent can extract additional learning from simulated experiences. However, the success of model-based methods hinges on the accuracy of the learned model, and inaccuracies may lead to suboptimal decisions.

In contrast, model-free methods learn value functions or policies directly from raw experiences, avoiding the need to model transition dynamics explicitly. Techniques such as Monte Carlo methods and Temporal Difference (TD) learning exemplify this approach. While model-free methods are generally simpler and more robust in complex or stochastic environments, they typically require more interactions with the environment to converge to an optimal solution. Recent developments—such as Dyna architectures—seek to combine these two paradigms, leveraging the sample efficiency of model-based approaches and the robustness of model-free methods.

## 2.4   Temporal Difference Learning Methods

Temporal difference (TD) learning stands out as a pivotal algorithmic breakthrough in reinforcement learning, merging ideas from dynamic programming and Monte Carlo methods. TD methods update value estimates based on the immediate reward and the estimated value of subsequent states—a process known as bootstrapping. For example, the TD(0) update rule is given by:

$$V(s_t) \leftarrow V(s_t) + \alpha \left[ r_{t+1} + \gamma V(s_{t+1}) - V(s_t) \right],$$

where $\alpha$ is the learning rate and the term in brackets is the TD error.

For control problems, TD learning has been extended to methods such as Sarsa (an on-policy method) and Q-learning (an off-policy method). Sarsa updates its action-value function as:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \Big[ r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t) \Big],$$

while Q-learning uses the update:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \Big[ r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t) \Big].$$

Both approaches have been shown to converge under suitable conditions, with Q-learning providing the flexibility to learn about optimal policies while following exploratory strategies.

TD learning further extends to TD($\lambda$) methods by introducing eligibility traces, which bridge the gap between one-step TD and Monte Carlo methods. The forward view of TD($\lambda$) defines the $\lambda$-return as:

$$G_t^\lambda = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_t^{(n)},$$

with the n-step return

$$G_t^{(n)} = \sum_{k=0}^{n-1} \gamma^k r_{t+k} + \gamma^n V(s_{t+n}).$$

The backward view implements this idea efficiently via eligibility traces, updating the value function as:

$$V(s) \leftarrow V(s) + \alpha \, \delta_t \, e(s),$$

where the eligibility trace $e(s)$ is updated by:

$$e(s) \leftarrow \gamma \lambda \, e(s) + \mathbf{1}\{s_t = s\}.$$

This approach has been instrumental in achieving improved convergence properties and handling both episodic and continuing tasks Sutton and Barto (2018); Huh and Mohapatra (2024).

## 2.5 Policy Optimization and Actor-Critic Methods

While value-based methods focus on estimating state or state-action values, policy-based methods directly optimize the policy. In Policy Gradient methods, the agent learns a parameterized policy and adjusts the parameters according to the gradient of the expected reward. The policy gradient theorem provides the foundation:

$$\nabla_\theta J(\theta) = \mathbb{E}_\pi \left[ \nabla_\theta \log \pi_\theta(a|s) \, Q^\pi(s, a) \right].$$

A common variant REINFORCE introduces a baseline $b(s)$ to reduce variance:

$$\nabla_\theta J(\theta) = \mathbb{E}_\pi \left[ \nabla_\theta \log \pi_\theta(a|s) \left( Q^\pi(s,a) - b(s) \right) \right].$$

Policy gradient methods are appealing because they directly optimize the policy without the need for an explicit value function. The reinforce algorithm is a Monte Carlo approach that uses the following update:

$$\theta_{t+1} = \theta_t + \alpha \, \nabla_\theta \log \pi_\theta(a_t|s_t) \, G_t,$$

where $G_t$ is the return following time $t$.

Variance reduction is critical in practice. One effective approach is the use of an *advantage function* $A^\pi(s,a) = Q^\pi(s,a) - V^\pi(s)$, which measures how much better an action is compared to the average. This leads to the following update:

$$\theta_{t+1} = \theta_t + \alpha \, \nabla_\theta \log \pi_\theta(a_t|s_t) \, A^\pi(s_t, a_t).$$

Actor-Critic methods combine these ideas by using a critic to estimate $V^\pi(s)$ or $Q^\pi(s,a)$ and an actor to update the policy. Recent advances in this area have led to algorithms with improved convergence properties and sample efficiency (Zeng et al., 2024a; Doan, 2021).

## 2.6 Epsilon-Greedy Policies and the Exploration-Exploitation Tradeoff

A fundamental challenge in reinforcement learning is managing the trade-off between exploration (seeking new information) and exploitation (using known information). A common strategy, the $\epsilon$-greedy policy, selects the best-known action with probability $1 - \epsilon$ while choosing a random action with probability $\epsilon$:

$$\pi(a|s) = \begin{cases} 1 - \epsilon + \frac{\epsilon}{|A|}, & \text{if } a = \arg\max_{a'} Q(s, a'), \\ \frac{\epsilon}{|A|}, & \text{otherwise.} \end{cases}$$

Beyond $\epsilon$-greedy, more sophisticated exploration techniques—such as softmax action selection, optimistic initialization, Upper Confidence Bounds (UCB), Thompson sampling, and intrinsic motivation strategies—have been developed. These methods address challenges like sparse rewards, deceptive local optima, and long-horizon planning, and they are crucial in complex environments where efficient exploration is essential for learning robust policiesSutton and Barto (2018).

## 2.7 Stochastic Approximation and Two-Time Scale Methods

Stochastic Approximation is a class of iterative methods for solving optimization problems when the underlying function is noisy. In RL, these methods are particularly useful when the environment's model is unknown or dynamic. Two-time scale methods enhance stochastic approximation by employing different

learning rates for the policy and value function updates, which can lead to improved stability and efficiency. A general stochastic approximation update is given by:

$$\theta_{t+1} = \theta_t + \alpha_t \Big( h(\theta_t) + M_{t+1} \Big),$$

where $\alpha_t$ is a diminishing step-size, $h(\theta_t)$ is the mean update, and $M_{t+1}$ represents noise.

In many RL algorithms (notably Actor-Critic), different parameters are updated at different rates. Suppose we have two sets of parameters: $\theta$ (e.g., for the actor) and $\omega$ (e.g., for the critic). Their updates are performed on different time scales:

$$\theta_{t+1} = \theta_t + \alpha_t \Big( h(\theta_t, \omega_t) + M_{t+1}^\theta \Big), \tag{1}$$

$$\omega_{t+1} = \omega_t + \beta_t \Big( g(\theta_t, \omega_t) + M_{t+1}^\omega \Big), \tag{2}$$

with the step-sizes satisfying

$$\lim_{t \to \infty} \frac{\alpha_t}{\beta_t} = 0.$$

This condition ensures that $\omega_t$ (the critic) converges much faster than $\theta_t$ (the actor), so that the actor's updates effectively see a near-converged critic. Two-time scale stochastic approximation has been rigorously analyzed using ordinary differential equation (ODE) techniques, and its general convergence properties have been established in recent works Doan and Romberg (2019); Zeng et al. (2024b); Borkar (2024).

These techniques are essential in ensuring the stability of complex RL algorithms where multiple interacting components are updated simultaneously, particularly in non-stationary or high-dimensional settings.

A major theoretical result in this area is provided by Borkar's Two Time-Scale Convergence Theorem. One common version of the theorem is stated as follows:

**Theorem 1** (Borkar's Two Time-Scale Convergence Theorem Borkar (2008)). *Consider the two-time scale stochastic approximation iterates given by*

$$\theta_{t+1} = \theta_t + \alpha_t \left( h(\theta_t, \omega_t) + M_{t+1}^\theta \right),$$
$$\omega_{t+1} = \omega_t + \beta_t \left( g(\theta_t, \omega_t) + M_{t+1}^\omega \right),$$

*where the step-sizes $\{\alpha_t\}$ and $\{\beta_t\}$ are positive, diminishing, and satisfy*

$$\lim_{t \to \infty} \frac{\alpha_t}{\beta_t} = 0.$$

*Under standard regularity conditions (including Lipschitz continuity of h and g, boundedness of the iterates, and appropriate martingale difference noise conditions), the following holds:*

1. *The fast iterate $\omega_t$ converges to a unique stable equilibrium $\omega^*(\theta)$ of the ODE*

$$\dot{\omega}(t) = g(\theta, \omega(t))$$

*for a fixed $\theta$.*

2. *The slow iterate $\theta_t$ converges to an invariant set of the ODE*

$$\dot{\theta}(t) = h\big(\theta(t), \omega^*(\theta(t))\big).$$

This theorem is pivotal in the analysis of Actor-Critic algorithms and other multi-time scale methods, ensuring that the slower parameter updates see a quasi-static environment as the faster updates converge rapidly. Such convergence guarantees are crucial in establishing the reliability and stability of learning algorithms in dynamic environments.

## 2.8 Summary

Reinforcement Learning is a dynamic and rapidly evolving field built upon the formalism of Markov Decision Processes, Dynamic Programming, and Temporal Difference learning. The integration of model-free methods, action-value approaches, and policy optimization techniques (including actor-critic and policy gradient methods) has enabled RL to tackle complex, high-dimensional decision-making tasks. Moreover, advanced tools such as TD($\lambda$) for efficient credit assignment and two-time scale stochastic approximation for stability in multi-parameter updates have been instrumental in the development of robust RL algorithms. As RL continues to find applications in robotics, game playing, healthcare, and autonomous systems, ongoing research remains focused on improving sample efficiency, convergence properties, and scalability.

# 3 Federated RL (FRL)

## 3.1 Introduction to FRL

Federated Reinforcement Learning (FRL) emerges as an innovative approach that integrates the privacy-preserving mechanisms of Federated Learning (FL) with the decision-making capabilities of Reinforcement Learning (RL) Qi et al. (2021). It represents a significant advancement in machine learning that addresses critical challenges in modern distributed systems where data privacy is paramount. According to the paper, FRL can be defined as "an integration of FL and RL under privacy protection, where several elements of RL can be presented in FL frameworks to deals with sequential decision-making tasks" Qi et al. (2021).

The fundamental architecture of FRL maintains the core principles of federated systems where multiple agents collaborate to build a shared model without directly exchanging their private data Qi et al. (2021). Instead, these agents

share model parameters or gradients, allowing the system to leverage collective experiences while maintaining data locality and privacy. This collaborative approach enables agents to benefit from the experiences of others without compromising sensitive information.

FRL algorithms are categorized into two main types based on the distribution characteristics of agents within the framework: Horizontal Federated Reinforcement Learning (HFRL) and Vertical Federated Reinforcement Learning (VFRL) Qi et al. (2021). This categorization mirrors the approach used in traditional federated learning, where similar distinctions are made based on data distribution patterns. The paper elaborates on these categories by providing specific case studies, including applications in autonomous driving and smart grid systems, to illustrate the practical implementation differences between HFRL and VFRL approaches Qi et al. (2021).

### 3.1.1   Core Components of FRL Systems

The FRL framework consists of distributed agents operating in potentially different environments, a coordination mechanism for model aggregation, and a secure communication protocol Qi et al. (2021). Each agent maintains its own local model trained on private experience data, while the aggregation process combines these models to create a globally improved policy that benefits all participants. This distributed architecture allows FRL to scale effectively across numerous devices or systems while maintaining robust privacy protections.

### 3.1.2   Relationship Between FL and RL

The integration of FL and RL creates a synergistic relationship that addresses limitations in both fields. FL contributes its distributed learning framework and privacy-preserving mechanisms, while RL provides methods for sequential decision-making in dynamic environments through interaction and feedback Qi et al. (2021).

Federated Learning emerged as a solution to privacy concerns in traditional machine learning approaches. It enables multiple parties to collaboratively train a model without sharing their raw data Qi et al. (2021). As described in the paper, FL is "a decentralized collaborative approach that allows multiple partners to train data respectively and build a shared model while maintaining privacy" Qi et al. (2021). The FL process typically involves local model training, secure aggregation of model updates, and redistribution of the improved model to participants.

Reinforcement Learning, conversely, focuses on how agents learn optimal behaviors through interactions with an environment. The paper defines RL as "one of the branches [of machine learning] that focuses on how individuals, i.e., agents, interact with their environment and maximize some portion of the cumulative reward" Qi et al. (2021). This learning process occurs through trial and error, with the agent developing policies that maximize expected rewards over time.

### 3.1.3 Integrating Reinforcement Learning Principles with Federated Architectures

Building upon the foundational concepts of Reinforcement Learning previously established, we now examine how these principles extend to federated environments. The integration of RL mechanisms with federated architectures necessitates reconceptualizing several core aspects of the traditional RL framework to accommodate distributed learning while preserving data privacy.

**Markov Decision Processes in FRL** In federated reinforcement learning, the traditional Markov Decision Process (MDP) framework requires adaptation to account for distributed agents operating under data isolation constraints. An FRL system can be formalized as a collection of local MDPs, where each agent $i$ maintains its own tuple $(\mathcal{S}_i, \mathcal{A}_i, \mathcal{P}_i, \mathcal{R}_i, \gamma)$ representing states, actions, transition probabilities, rewards, and discount factor, respectively Qi et al. (2021).

The key distinction emerges in how these local MDPs interact within the federated framework. McMahan et al. (2017), who introduced the seminal FedAvg algorithm, provide the foundation for parameter aggregation in federated settings. Building upon and extending these foundations, Jin et al. (2022a) propose that policy interactions in FRL can be modeled as:

$$\pi_{\text{global}} = \mathcal{F}(\{\pi_i\}_{i=1}^N, \{w_i\}_{i=1}^N)$$

where $\pi_i$ represents the policy of agent $i$, $w_i$ is the weight assigned to agent $i$ (typically proportional to the size of its local dataset), and $\mathcal{F}$ is the federation function that aggregates local policies into a global policy $\pi_{\text{global}}$ using weighted averaging similar to FedAvg.

**Value Function Decomposition in FRL** Following policy aggregation, value function decomposition becomes essential in federated contexts. Rather than maintaining a single centralized value function, FRL systems operate with a distributed set of value functions that must be coherently aggregated. The work Zhuo et al. (2019) formulates this relationship through a weighted averaging approach:

$$V_{\text{global}}^{\pi}(s) = \sum_{i=1}^N \frac{n_i}{n_{\text{total}}} V_i^{\pi}(s)$$

where $V_i^{\pi}(s)$ represents the value function of agent $i$ under policy $\pi$ for state $s$, $n_i$ represents the number of samples from agent $i$, and $n_{\text{total}} = \sum_{i=1}^N n_i$ is the total number of samples across all agents. This approach ensures that agents with more extensive experience have proportionally greater influence on the global value function, similar to the FedAvg mechanism proposed by McMahan et al. (2017).

**Federated Policy Optimization**   The policy optimization process in FRL requires balancing local improvements with global convergence. Chen et al. Chen et al. (2021) propose a federated policy gradient approach that modifies the traditional policy gradient theorem to accommodate federated constraints:

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}_{\pi_\theta}\left[\sum_{i=1}^{N} \alpha_i \nabla_\theta \log \pi_\theta(a|s) Q_i^\pi(s,a)\right] \tag{3}$$

where $\alpha_i$ represents the contribution factor of agent $i$ (typically proportional to dataset size), and $Q_i^\pi(s,a)$ is the local action-value function. This approach enables agents to optimize policies based on their local experiences while contributing to a globally improved policy. Importantly, Chen et al. provide theoretical convergence guarantees for this method under certain assumptions, although convergence rates may degrade with highly non-IID data distributions across agents.

**Privacy-Preserving Experience Replay**   Experience replay, a technique that enhances sample efficiency in RL by reusing past experiences, faces unique challenges in federated settings due to privacy constraints. Traditional experience sharing is prohibited in FRL, necessitating alternative approaches. Truex et al. (2019) introduce hybrid differential privacy mechanisms applicable to federated learning, which Wei et al. (2020) extend to reinforcement learning contexts:

$$D_{\text{fed}} = \psi(\{D_i\}_{i=1}^{N}, \epsilon, \delta) \tag{4}$$

where $D_i$ represents the local experience buffer of agent $i$, and $\psi$ is a privacy-preserving aggregation function that uses differential privacy with privacy budget parameters $\epsilon$ and $\delta$. This approach enables knowledge transfer without raw data exchange by adding calibrated noise to statistical summaries rather than sharing raw experiences, thereby maintaining $(\epsilon, \delta)$-differential privacy guarantees.

**Federated Q-Learning**   Building on experience replay mechanisms, Q-learning adapts to federated environments through modifications to its update rule. Nadiger et al. (2019) propose a federated Q-learning approach with the following update mechanism:

$$Q_{\text{global}}(s,a) \leftarrow (1-\beta)Q_{\text{global}}(s,a) + \beta \sum_{i=1}^{N} \frac{n_i}{n_{\text{total}}} Q_i(s,a) \tag{5}$$

where $\beta$ is the global learning rate, $n_i$ is the number of samples contributed by agent $i$, and $n_{\text{total}}$ is the total number of samples across all agents. This approach weights the contributions of agents based on their local dataset sizes. Notably, Nadiger et al. acknowledge an important challenge: when agents have

conflicting Q-values for the same state-action pairs due to environment differences, simple averaging may lead to suboptimal policies. They propose a confidence-weighted aggregation as a partial solution, but recognize this as an ongoing research challenge in heterogeneous FRL settings.

### 3.1.4 Mapping FL and RL Components

FRL leverages FL's decentralized learning framework while addressing sequential decision-making tasks inherent in RL. The components of FL are mapped to RL elements as follows:

- The environment in Reinforcement Learning (RL) corresponds to the sample space in Federated Learning (FL).

- States observed by agents correspond to features in traditional FL.

- Actions selected by agents correspond to labels or predictions in FL models.

This mapping ensures that FRL retains the strengths of both FL and RL while addressing their limitations Qi et al. (2021).

By integrating these RL principles into FL frameworks, FRL enables distributed agents to collaboratively solve complex sequential decision-making tasks while preserving privacy. This synergy between FL and RL represents a significant advancement in machine learning for modern distributed systems.

## 3.2 FRL Architecture and Categories

According to the distribution characteristics of agents within the framework, FRL architectures can be systematically classified into two fundamental categories:

1. **Horizontal Federated Reinforcement Learning (HFRL)**: Addresses scenarios where agents share the same state-action spaces but possess different experiences or data distributions.

2. **Vertical Federated Reinforcement Learning (VFRL)**: Manages scenarios where different agents observe different features or aspects of the environment, requiring integration of partial observations.

This categorization parallels the classification in Federated Learning between Horizontal Federated Learning (HFL) and Vertical Federated Learning (VFL), but adapts to the sequential decision-making context of reinforcement learning.

### 3.2.1 Horizontal Federated Reinforcement Learning (HFRL)

**Theoretical Foundation** HFRL addresses scenarios where multiple agents operate in similar environments with consistent state-action spaces but different experiences or data distributions. The autonomous driving case exemplifies

this architecture, where vehicles navigate similar action and state spaces but encounter different environmental conditions.

In HFRL, the agents share the same state and action spaces, but may experience different transition dynamics and reward functions based on their unique environment interactions:

$$\forall i, j : S_i = S_j \text{ and } A_i = A_j \tag{6}$$

But potentially different transition and reward functions:

$$\exists i, j : P_i \neq P_j \text{ or } R_i \neq R_j \tag{7}$$

**Mathematical Formulation** For an HFRL system with $N$ agents, each agent $i$ operates within its own MDP:

$$\text{MDP}_i = (S, A, P_i, R_i, \gamma) \tag{8}$$

Each agent aims to maximize its expected return based on a policy $\pi_i$:

$$J(\pi_i) = \mathbb{E}_{\tau \sim \pi_i} \left[ \sum_{t=0}^{\infty} \gamma^t R_i(s_t, a_t, s_{t+1}) \right] \tag{9}$$

Where $\tau$ represents trajectories generated by following policy $\pi_i$.

**Value Function Aggregation** Using Q-learning as an example implementation, each agent maintains a local Q-function $Q_i(s, a)$ that is updated through local experiences:

$$Q_i(s_t, a_t) \leftarrow Q_i(s_t, a_t) + \alpha \left[ r_t + \gamma \max_a Q_i(s_{t+1}, a) - Q_i(s_t, a_t) \right] \tag{10}$$

The federated aggregation periodically combines these local functions into a global Q-function:

$$Q_{\text{global}}(s, a) = \sum_{i=1}^{N} \frac{n_i}{n} Q_i(s, a) \tag{11}$$

This aggregated function is then distributed back to all agents for continued local learning.

**Policy-Based Formulation** For policy-based methods, each agent maintains a parameterized policy $\pi_i(a|s; \theta_i)$ where $\theta_i$ represents policy parameters. The local objective function for agent $i$ is:

$$J_i(\theta_i) = \mathbb{E}_{s \sim \rho^{\pi_i}, a \sim \pi_i(\cdot|s; \theta_i)} \left[ Q^{\pi_i}(s, a) \right] \tag{12}$$

Where $\rho^{\pi_i}$ denotes the state distribution under policy $\pi_i$. The local policy parameters are updated using gradient ascent:

$$\theta_i \leftarrow \theta_i + \beta \nabla_{\theta_i} J_i(\theta_i) \tag{13}$$

The federation process aggregates these parameters:

$$\theta_{\text{global}} = \sum_{i=1}^{N} \frac{n_i}{n} \theta_i \tag{14}$$

This approach enables collaborative policy improvement while preserving the privacy of local experiences.
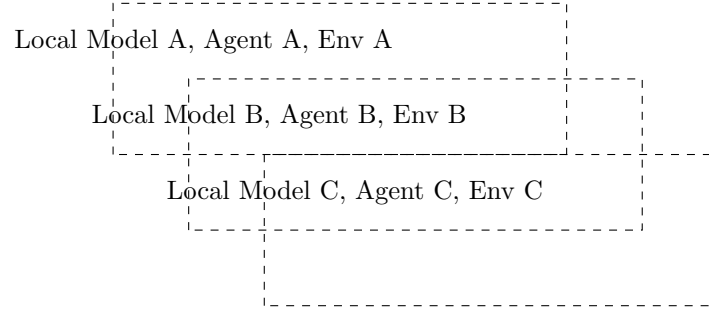
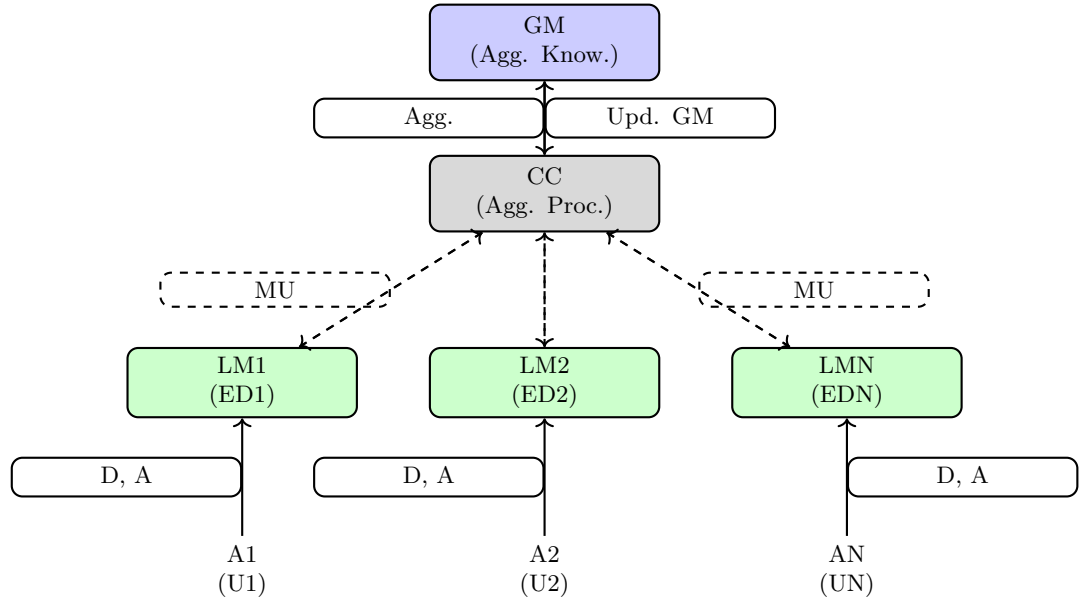Figure 1: Horizontal Federated Reinforcement Learning



Figure 2: An example of horizontal federated reinforcement learning (HFRL) architecture.

| Abbr. | Meaning |
|---|---|
| GM | Global Model |
| CC | Central Coordinator |
| LM | Local Model |
| ED | Edge Device |
| A | Agent |
| U | User |
| D | Data |
| MU | Model Update |
| Agg. | Aggregation |
| Upd. GM | Updated Global Model |

### 3.2.2 Vertical Federated Reinforcement Learning (VFRL)

**Theoretical Foundation**  VFRL addresses scenarios where different agents observe different features or aspects of the environment, requiring integration of these partial observations. The smart grid case exemplifies this architecture, where different components of the grid observe different aspects of the overall system state.

In VFRL, the environment state $s$ is partitioned across agents, with each agent $i$ observing only a subset $s^i$ of the complete state:

$$s = \{s^1, s^2, \ldots, s^N\} \tag{15}$$

Similarly, the action space may be partitioned:

$$a = \{a^1, a^2, \ldots, a^N\} \tag{16}$$

This can be formally characterized as:

$$\forall i \neq j : S^i \cap S^j = \emptyset \text{ or } A^i \cap A^j = \emptyset \tag{17}$$

This formulation represents the extreme case of complete partitioning, though partial overlaps may exist in practical implementations.

**Mathematical Formulation**  The complete MDP for a VFRL system can be represented as:

$$\text{MDP} = (S = S^1 \times S^2 \times \ldots \times S^N, A = A^1 \times A^2 \times \ldots \times A^N, P, R, \gamma) \tag{18}$$

Each agent maintains a local value function or policy based on its partial observation. For value-based methods, agent $i$ maintains $Q_i(s^i, a)$, which estimates expected returns given its observed state component.

**Feature-Based Integration**  A sophisticated VFRL approach involves agents exchanging feature representations rather than direct value functions. Each agent $i$ maintains an encoder network $E_i(s^i)$ that maps its partial observation to a feature vector $\phi_i$. These feature vectors are exchanged and combined:

$$\phi = g(\phi_1, \phi_2, \ldots, \phi_N) \tag{19}$$

Where $g$ is a feature aggregation function. This combined representation is then used to estimate the value function or policy:

$$Q(s, a) = h(\phi, a) \tag{20}$$

The parameters of encoders $E_i$, aggregator $g$, and value estimator $h$ are jointly optimized through federation while preserving privacy of raw observations.

**Combined Value Function Approximation**  The federated process in VFRL involves combining partial value functions to approximate the global value function:

$$Q(s, a) \approx f(Q_1(s^1, a), Q_2(s^2, a), \ldots, Q_N(s^N, a)) \tag{21}$$

Where $f$ represents an aggregation function that combines the partial value estimates. This function could be a simple average, a weighted sum, or a more complex neural network that learns to optimally combine the partial estimates.
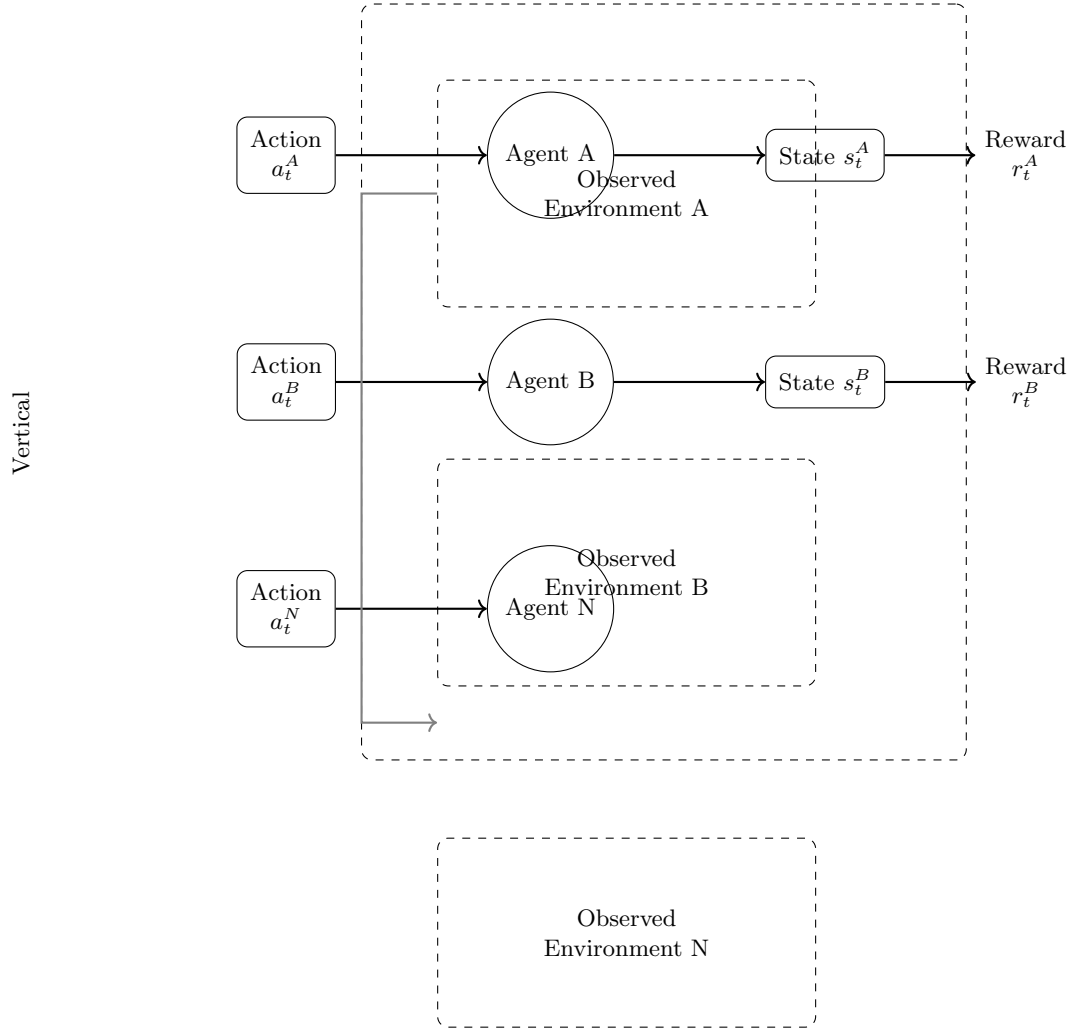
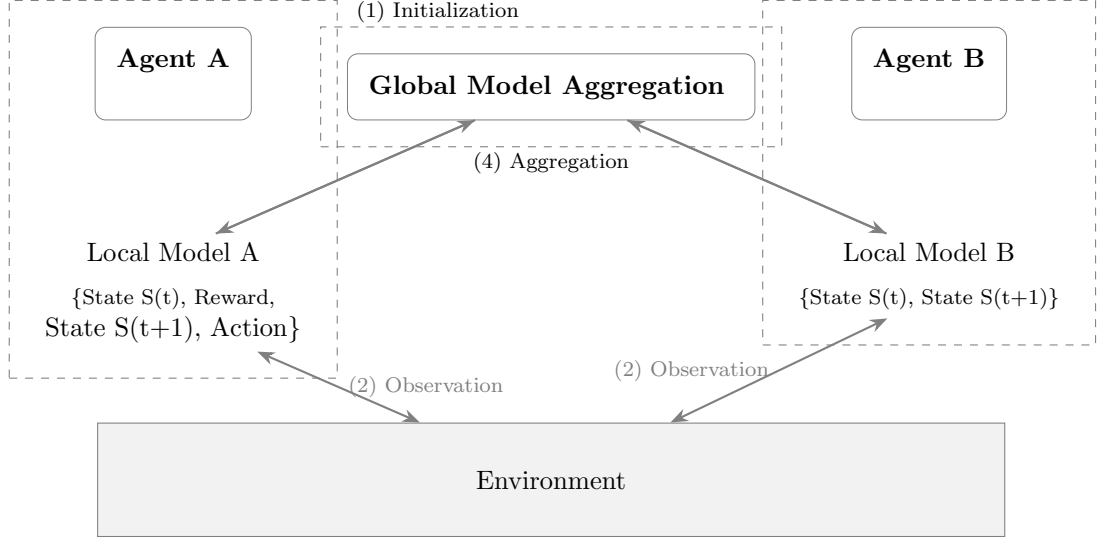Figure 3: Vertical Federated Reinforcement Learning

Figure 4: An example of vertical federated reinforcement learning (VFRL) architecture

## 3.3 Advantages of FRL

FRL offers several significant advantages over traditional RL approaches, particularly in distributed and privacy-sensitive applications. These advantages directly address critical challenges in implementing RL in real-world scenarios.

### 3.3.1 Privacy Preservation with Effective Information Exchange

One of the primary advantages of FRL is its ability to facilitate information exchange between agents while protecting privacy. The paper states that "FL can not only complete information exchange while avoiding privacy disclosure, but also adapt various agents to their different environments" Qi et al. (2021). This privacy-preserving mechanism is essential in applications where agent data may contain sensitive information that should not be shared directly with other participants or even with a central server.

### 3.3.2 Bridging the Simulation-Reality Gap

A significant challenge in traditional RL is the disparity between simulated training environments and real-world deployment conditions. FRL offers a solution to this problem as it "can aggregate information from both environments and thus bridge the gap between them" Qi et al. (2021). By enabling collaboration between agents operating in both simulated and real environments, FRL creates more robust and adaptable policies that perform well in actual deployment scenarios.

### 3.3.3 Addressing Sample Efficiency and Large State-Action Spaces

FRL provides solutions to core RL challenges related to sample efficiency and large state-action spaces. The paper notes that "in the case of large action space and state space, the performance of agents is vulnerable to collected samples since it is nearly impossible to explore all sampling spaces" Qi et al. (2021). By enabling agents to share learning experiences through model parameters rather than raw data, FRL significantly improves sample efficiency and accelerates the learning process.

### 3.3.4 Integration of Partial Observations

In many practical scenarios, individual agents have access to only partial observations of the environment. FRL addresses this limitation by enabling the integration of these partial observations through secure aggregation. The paper highlights that "in some cases, only partial features can be observed by each agent in RL... FL makes it possible to integrate this information through aggregation" Qi et al. (2021). This capability allows for more comprehensive decision-making based on collective information without requiring direct data sharing.

## 3.4 Federated RL Communication Structures

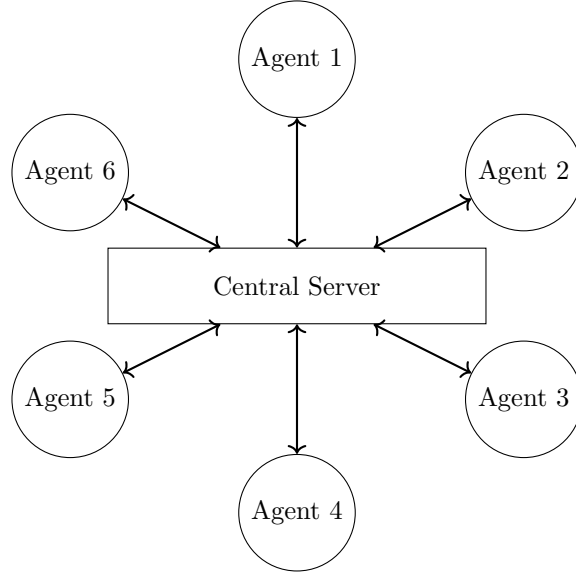### 3.4.1 Star Communication (Centralized Aggregation)

Star communication follows a client-server architecture, where a central aggregator (server) collects model updates from multiple agents, processes them, and distributes the updated model back to all agents Jin et al. (2022b); Lyu et al. (2021).

In this architecture, each agent independently trains a local model using its unique dataset. Instead of directly exchanging model parameters with other agents, they send updates to a central server. The server performs model aggregation using techniques such as FedAvg Zhuo et al. (2019); Zhang et al. (2023), where a weighted averaging method is applied to combine updates from all agents. The newly aggregated model is then redistributed to all agents, ensuring that each participant benefits from the collective learning process. This iterative process continues until the global model converges.

One of the main advantages of this approach is its scalability. The central server efficiently manages communication between agents, reducing network congestion. Additionally, since agents do not communicate with each other directly, privacy is preserved because only model updates are shared rather than raw data Lyu et al. (2021). This setup ensures that all agents follow a unified learning trajectory, leading to consistent policy updates and stability in learning Zhang et al. (2023).

However, this approach has its limitations. The central server becomes a single point of failure, meaning if it crashes or becomes compromised, the entire learning process halts Lyu et al. (2021). Furthermore, as the number of

agents increases, the server can experience bottlenecks, slowing down aggregation and increasing latency Zhang et al. (2023). Lastly, since all agents receive the same aggregated model, they may struggle to adapt to highly heterogeneous environments where local variations are significant.



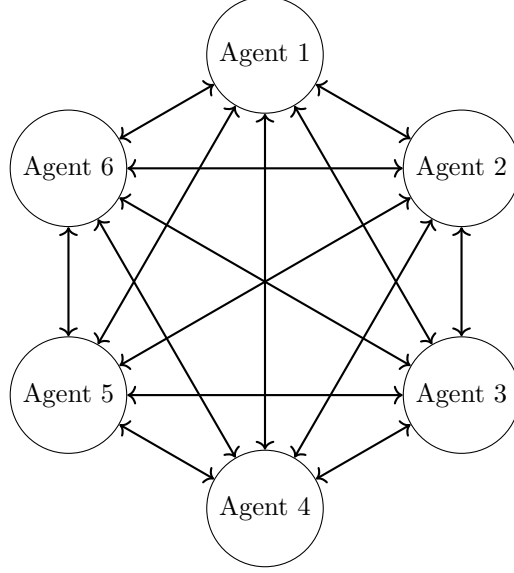### 3.4.2 All-to-All Communication (Decentralized Aggregation)

All-to-all communication eliminates the need for a central aggregator. Instead, agents exchange model updates directly with their peers in a distributed network topology Zhuo et al. (2019); Lyu et al. (2021).

Each agent trains a local model using its dataset and periodically communicates with selected peers, rather than relying on a central server. Model updates are shared among neighboring agents and aggregated using decentralized consensus-based techniques such as decentralized averaging Jin et al. (2022b); Zhang et al. (2023). The process continues iteratively, with each agent incorporating the updates received from its peers to refine its model until convergence is reached.

One of the primary benefits of all-to-all communication is its robustness. Since there is no central authority, the system can continue operating even if some agents drop out, making it highly fault-tolerant Lyu et al. (2021). This approach is particularly beneficial in dynamic environments where agents operate under different conditions. Additionally, privacy is enhanced because model updates are exchanged only between direct peers rather than a central server Zhuo et al. (2019). This setup is well-suited for applications such as distributed sensor networks and multi-agent robotic systems.

However, this method also presents challenges. Communication overhead is significantly higher compared to star communication because each agent must

exchange updates with multiple peers, leading to increased bandwidth consumption Zhang et al. (2023). Furthermore, maintaining synchronization between all agents is difficult due to differences in local processing speeds and network conditions Lyu et al. (2021). Additionally, achieving convergence in decentralized networks can be slower compared to centralized aggregation, as updates propagate gradually across the network rather than being centrally applied at once.



## 3.5 Federated Reinforcement Learning Algorithms

The following algorithms are based on the Federated Reinforcement Learning (FedRL) framework introduced by Jin et al. in their paper "Federated Reinforcement Learning with Environment Heterogeneity" Jin et al. (2022b). This framework addresses scenarios where multiple agents operate in environments with different state transition dynamics while sharing the same state space, action space, and reward function. The main challenge is collaborative learning without sharing raw trajectory data to preserve privacy. The objective is formalized as optimizing a policy that maximizes the average expected return across all environments:

$$g_{d_0}(\pi) = \frac{1}{n} \sum_{i=1}^{n} \mathbb{E}\left[\sum_{t=1}^{\infty} \gamma^t R(s_t, a_t) | s_0 \sim d_0, a_t \sim \pi(\cdot|s_t), s_{t+1} \sim P_i(\cdot|s_t, a_t)\right]$$
(22)

Where $d_0$ represents the initial state distribution, $\pi$ is the policy being optimized, and $P_i$ represents the transition dynamics of environment $i$. Jin et al. prove that no single policy can be simultaneously optimal across all het-

erogeneous environments, as the optimal policy depends on the initial state distribution.

### 3.5.1 QAvg Algorithm: Federated Q-Learning

QAvg is a federated version of Q-learning, where agents train local Q-values and periodically aggregate them at a server Jin et al. (2022b). Each agent updates its local Q-function using its environment's dynamics. The update combines a standard Q-learning update with a momentum term $(1 - \eta_t)$ that helps stabilize learning across heterogeneous environments. After several local updates, the server averages the Q-tables from all agents and redistributes this averaged Q-table back to each agent. Only the Q-values are communicated, preserving the privacy of each agent's experiences.

---

**Algorithm 1** QAvg Algorithm (Federated Q-Learning)Jin et al. (2022b)

---

1: **Initialize:** Each agent $k$ initializes $Q_0^k(s, a)$ arbitrarily for all states $s$ and actions $a$
2: **for** each round $t = 0, 1, 2, \ldots$ **do**
3:     **Local Updates:** Each agent $k$ performs multiple local updates:
4:     $Q_{t+1}^k(s, a) \leftarrow (1 - \eta_t) \cdot Q_t^k(s, a) + \eta_t \cdot \left[ R(s, a) + \gamma \sum_{s'} P_k(s'|s, a) \max_{a'} Q_t^k(s', a') \right]$
5:     **Global Aggregation:** The server performs:
6:     $\bar{Q}_t(s, a) \leftarrow \frac{1}{n} \sum_{k=1}^{n} Q_t^k(s, a)$, for all $s, a$
7:     **Distribution:** For all agents $k$ and all state-action pairs:
8:     $Q_t^k(s, a) \leftarrow \bar{Q}_t(s, a)$
9: **end for**

---

Jin et al. theoretically prove that QAvg converges to a suboptimal solution where the degree of suboptimality depends on the heterogeneity of the environments. The algorithm is most effective when the state transition functions across environments are similar.

### 3.5.2 PAvg Algorithm: Federated Policy Gradient

PAvg extends policy gradient methods to the federated setting, allowing agents to update local policies while aggregating globally Jin et al. (2022b). Each agent independently performs policy gradient updates based on its local environment. After computing the gradient, the policy is projected back onto the action probability simplex to ensure it remains a valid probability distribution. Following several local updates, the server aggregates the policies from all agents and distributes the averaged policy back. Similar to QAvg, only the policy functions are shared, not the raw trajectories.

**Algorithm 2** PAvg Algorithm (Federated Policy Gradient)Jin et al. (2022b)

---

1: **Initialize:** Each agent $k$ initializes policy $\pi_0^k(a|s)$ for all states $s$ and actions $a$
2: **for** each round $t = 0, 1, 2, \ldots$ **do**
3:     **Local Updates:** Each agent $k$ performs multiple local updates:
4:     $\tilde{\pi}_{t+1}^k(a|s) \leftarrow \pi_t^k(a|s) + \frac{\partial g_{d_0,k}(\pi_t^k)}{\partial \pi(a|s)}$, for all $s, a$
5:     $\pi_{t+1}^k(\cdot|s) \leftarrow \text{Proj}_{\Delta(A)}(\tilde{\pi}_{t+1}^k(\cdot|s))$, for all $s$
6:     **Global Aggregation:** The server performs:
7:     $\bar{\pi}_t(a|s) \leftarrow \frac{1}{n} \sum_{k=1}^n \pi_t^k(a|s)$, for all $s, a$
8:     **Distribution:** For all agents $k$ and all state-action pairs:
9:     $\pi_t^k(a|s) \leftarrow \bar{\pi}_t(a|s)$
10: **end for**

---

The theoretical analysis by Jin et al. shows that PAvg also converges to a suboptimal solution with suboptimality dependent on environment heterogeneity. The algorithm provides a direct way to optimize policies when the optimal policy structure is known or when policy-based methods are preferred.

### 3.5.3 DQNAvg: Federated Deep Q-Networks

DQNAvg extends QAvg to deep reinforcement learning scenarios where the state space is too large for tabular representations Jin et al. (2022b). Each agent maintains a local Q-network that is updated using standard DQN techniques including experience replay and target networks. Instead of sharing Q-tables as in QAvg, agents share the parameters of their neural networks. The server periodically aggregates these parameters and redistributes them to all agents.

---

**Algorithm 3** DQNAvg (Federated Deep Q-Networks)Jin et al. (2022b)

---

1: **Initialize:** Each agent $k$ initializes Q-network parameters $\theta^k$
2: **for** each round $t = 0, 1, 2, \ldots$ **do**
3:     **Local Updates:** Each agent $k$ performs multiple DQN updates:
4:     Collect experiences $(s, a, r, s')$ by interacting with environment $k$
5:     Update Q-network using standard DQN techniques:
6:     $\theta^k \leftarrow \theta^k - \alpha \nabla_{\theta^k} \left( r + \gamma \max_{a'} Q_{\theta^k}(s', a') - Q_{\theta^k}(s, a) \right)^2$
7:     **Global Aggregation:** The server averages network parameters:
8:     $\theta_{\text{global}} \leftarrow \frac{1}{n} \sum_{k=1}^n \theta^k$
9:     **Distribution:** For all agents $k$:
10:     $\theta^k \leftarrow \theta_{\text{global}}$
11: **end for**

---

This approach enables FedRL to scale to more complex environments with large or continuous state spaces that cannot be effectively represented in tabular form. The neural network architecture allows for function approximation and better generalization across states.

### 3.5.4 DDPGAvg: Federated Deep Deterministic Policy Gradient

DDPGAvg extends PAvg to continuous action spaces using the Deep Deterministic Policy Gradient algorithm Jin et al. (2022b). Each agent maintains local actor and critic networks that are updated using standard DDPG techniques. After several local updates, the server aggregates the parameters of both networks across all agents and redistributes them.

---

**Algorithm 4** DDPGAvg (Federated Deep Deterministic Policy Gradient)Jin et al. (2022b)

---

1: **Initialize:** Each agent $k$ initializes actor network $\theta^{\pi^k}$ and critic network $\theta^{Q^k}$
2: **for** each round $t = 0, 1, 2, \ldots$ **do**
3:     **Local Updates:** Each agent $k$ performs multiple DDPG updates:
4:     Collect experiences $(s, a, r, s')$ by interacting with environment $k$
5:     Update critic network:
6:     $\theta^{Q^k} \leftarrow \theta^{Q^k} - \alpha_Q \nabla_{\theta^{Q^k}} \left( r + \gamma Q_{\theta^{Q^k}}(s', \pi_{\theta^{\pi^k}}(s')) - Q_{\theta^{Q^k}}(s, a) \right)^2$
7:     Update actor network:
8:     $\theta^{\pi^k} \leftarrow \theta^{\pi^k} + \alpha_\pi \nabla_{\theta^{\pi^k}} Q_{\theta^{Q^k}}(s, \pi_{\theta^{\pi^k}}(s))$
9:     **Global Aggregation:** The server averages network parameters:
10:     $\theta^{\pi}_{\text{global}} \leftarrow \frac{1}{n} \sum_{k=1}^{n} \theta^{\pi^k}$
11:     $\theta^{Q}_{\text{global}} \leftarrow \frac{1}{n} \sum_{k=1}^{n} \theta^{Q^k}$
12:     **Distribution:** For all agents $k$:
13:     $\theta^{\pi^k} \leftarrow \theta^{\pi}_{\text{global}}$
14:     $\theta^{Q^k} \leftarrow \theta^{Q}_{\text{global}}$
15: **end for**

---

This algorithm is particularly useful for continuous control tasks where the action space is continuous rather than discrete. The actor network directly outputs actions rather than action probabilities, and the critic network provides value estimates to guide the actor's updates.

### 3.5.5 Personalized Federated RL (PerDQNAvg & PerDDPGAvg)

Recognizing the limitations of a one-size-fits-all approach, Jin et al. propose personalized versions of their algorithms through environment embeddings Jin et al. (2022b). During training, agents share all network parameters except their embedding layers, allowing them to learn a common policy structure while retaining environment-specific information.

---

**Algorithm 5** Personalized Federated RL (PerDQNAvg & PerDDPGAvg)Jin et al. (2022b)

---

1: **Initialize:** Each agent $k$ initializes network parameters $\theta^k$ and environment embedding $e^k$
2: **for** each round $t = 0, 1, 2, \ldots$ **do**
3:     **Local Updates:** Each agent $k$ performs multiple updates:
4:     Collect experiences $(s, a, r, s')$ by interacting with environment $k$
5:     Update network using standard techniques, conditioning on environment embedding $e^k$:
6:     $\pi^k(a|s, e^k) = \pi_\theta(a|s, e^k)$
7:     Update embedding vector using gradient descent:
8:     $e^k \leftarrow e^k - \eta\nabla_e J(e)$
9:     **Global Aggregation:** The server averages all parameters except embeddings:
10:     $\theta_{\text{global}} \leftarrow \frac{1}{n}\sum_{k=1}^{n}\theta^k$ (excluding embedding layers)
11:     **Distribution:** For all agents $k$:
12:     $\theta^k \leftarrow \theta_{\text{global}}$ (keeping original embedding $e^k$)
13: **end for**

---

This approach enables better performance in individual environments while still benefiting from collaborative learning. Additionally, it facilitates generalization to new environments by only requiring adjustment of the embedding layer rather than retraining the entire model.

## 3.6 Theoretical Results

From Jin et al. (2022b), the results pertain to the behavior of the proposed algorithms QAvg and PAvg in federated reinforcement learning (FedRL) under heterogeneous environments.

**Theorem 1: Dependence of Optimal Policy on Initial State Distribution**

**Theorem 1.** There exists a task of FedRL with the following properties. Assume that $\pi^\star \in \arg\max_\pi g_{d_0}(\pi)$. There exists another initial state distribution $d_0'$ and another policy $\tilde{\pi}$ such that:

$$g_{d_0}(\tilde{\pi}) < g_{d_0}(\pi^\star), \quad \text{but} \quad g_{d_0'}(\tilde{\pi}) > g_{d_0'}(\pi^\star).$$

*Explanation:* This theorem demonstrates that in FedRL with heterogeneous environments, there does not exist a single optimal policy $\pi^\star$ that is universally optimal across all initial state distributions. The optimal policy depends on the specific initial state distribution $d_0$.

**Theorem 2: Convergence to Suboptimal Solutions**

**Theorem 2.** Both QAvg and PAvg converge to suboptimal solutions, and the suboptimality is affected by the degree of environment heterogeneity in FedRL.

*Explanation:* This result establishes that while QAvg and PAvg algorithms are guaranteed to converge, they converge to solutions that are suboptimal compared to environment-specific optimal policies. The extent of suboptimality is directly influenced by the heterogeneity of the environments.

## 3.7 Summary

| Method | Type of RL Algorithm | Theoretical Guarantee | Key Properties | Comm Structure | Applicability |
|---|---|---|---|---|---|
| QAvg | Value-based (Tabular Q-Learning) | Converges to a suboptimal solution; suboptimality depends on environment heterogeneity. | Aggregates Q-tables; preserves privacy by sharing only Q-values. | Star (Centralized Aggregation) | Discrete state-action spaces. |
| PAvg | Policy-based (Policy Gradient) | Converges to a suboptimal solution; suboptimality depends on environment heterogeneity. | Aggregates policy parameters; preserves privacy by sharing only policies. | Star (Centralized Aggregation) | Discrete or continuous action spaces. |
| DQNAvg | Value-based (Deep Q-Networks) | Extends QAvg to large/continuous state spaces; suboptimality due to heterogeneity. | Aggregates neural network parameters; uses experience replay and target networks. | Star (Centralized Aggregation) | Large or continuous state spaces. |
| DDPGAvg | Policy-based (DDPG) | Extends PAvg to continuous action spaces; suboptimality due to heterogeneity. | Aggregates actor-critic network parameters; suitable for continuous control tasks. | Star (Centralized Aggregation) | Continuous action spaces. |

| Method | Type of RL Algorithm | Theoretical Guarantee | Key Properties | Comm Structure | Applicability |
|---|---|---|---|---|---|
| PDQNAvg | Personalized (DQN) | No explicit guarantees; improves performance in heterogeneous environments. | Shares all network parameters except embedding layers; retains environment-specific adaptations. | Star (Centralized Aggregation) | Environments requiring personalization. |
| PDDPGAvg | Personalized (DDPG) | No explicit guarantees; improves performance in heterogeneous environments. | Similar to PerDQNAvg but for continuous actions; combines shared policy structure with local embeddings. | Star (Centralized Aggregation) | Continuous actions with personalization. |
| HFRL | General Framework | Depends on underlying algorithm (e.g., QAvg, PAvg). | Agents share state-action spaces but have different experiences; aggregates Q-values or policies. | Star or All-to-All | Homogeneous state-action spaces. |
| VFRL | General Framework | Depends on underlying algorithm (e.g., DQNAvg, DDPGAvg). | Agents observe different features; integrates partial observations via feature or value aggregation. | Star or All-to-All | Heterogeneous or partitioned environments. |

## 3.8   Open Problems and Future Work

Federated Reinforcement Learning (FedRL) with environment heterogeneity represents a burgeoning field that combines the challenges of distributed learning with the complexities of diverse environmental dynamics. While recent advancements, such as the QAvg and PAvg algorithms proposed by Jin et al. Jin et al. (2022b), have laid foundational frameworks for collaborative policy optimization across heterogeneous environments, significant open problems remain. This section synthesizes unresolved challenges and outlines promising directions for future research, drawing insights from theoretical analyses, algo-

rithmic limitations, and practical considerations highlighted in the literature.

## 3.9 Theoretical Foundations and Limitations

The convergence guarantees of QAvg and PAvg algorithms reveal a fundamental dependency on environment heterogeneity in federated reinforcement learning systems. Jin et al. demonstrate that both approaches inevitably converge to suboptimal policies, with performance gaps directly proportional to the divergence between state transition dynamics across individual environments Jin et al. (2022b). This finding quantifies the inherent trade-off between collaborative learning and environmental heterogeneity, raising several critical questions about potential improvements to the current paradigm. Current aggregation mechanisms employ naive averaging of local Q-tables or policies, which potentially dilutes environment-specific nuances that could be critical for optimal performance. Future research might explore more sophisticated weighted averaging schemes where aggregation weights dynamically adjust based on measured pairwise environment similarities or task affinities, potentially mitigating the heterogeneity-induced suboptimality that plagues current approaches.

The influence of initial state distributions further complicates the FedRL landscape. As demonstrated in Theorem 1 of Jin et al., the optimal policy in FedRL exhibits dependency on the initial state distribution $d_0$ Jin et al. (2022b). This dependence significantly complicates policy generalization across environments with divergent starting conditions, necessitating more robust analytical frameworks that can account for these distribution shifts. The theoretical framework introduces an important construct-the imaginary environment $M_I$ with averaged transition dynamics $\bar{P}$-which serves as an analytical proxy for convergence analysis. However, this construct introduces its own limitations, as $M_I$ may not correspond to any actual real-world environment, thereby constraining the interpretability and practical applicability of the theoretical guarantees. Developing more formalized relationships between $M_I$ and the original environments could potentially yield tighter performance bounds and more actionable insights for practitioners implementing these systems.

The existing convergence analyses predominantly assume tabular representations of Q-functions and policies. This assumption diverges significantly from modern reinforcement learning applications, which overwhelmingly employ function approximators such as deep neural networks. While Jin et al. propose extended approaches like DQNAvg and DDPGAvg for non-tabular settings, the theoretical properties of these extensions remain largely uncharacterized Jin et al. (2022b). This gap raises fundamental questions about convergence conditions in heterogeneous settings when using neural network approximations. The relationship between network overparameterization and the critical trade-off between generalization and personalization remains unexplored, as does the potential for approximation errors introduced by neural architectures to amplify the suboptimality issues already present due to environment heterogeneity.

## 3.10 Personalization Approaches and Challenges

The personalization heuristic proposed by Jin et al.-embedding environments into low-dimensional vectors-demonstrates empirical performance improvements but currently lacks comprehensive theoretical justification Jin et al. (2022b). Bridging this theoretical gap requires formal characterization of the embedding spaces to understand what geometric or topological properties these environment embeddings should satisfy to effectively capture transition dynamics. These embeddings might be productively related to existing concepts in reinforcement learning literature, such as bisimulation metrics, providing a firmer foundation for the approach. Additionally, while the global policy converges to a suboptimal solution, the complex interaction between personalized embeddings and global network parameters remains unanalyzed. A two-timescale analysis framework, where embeddings evolve faster than global weights, could provide valuable insights into this relationship and guide more effective algorithm design.

The heuristic allows for fine-tuning embeddings when encountering new environments, but theoretical limits on this adaptability remain undefined. Developing information-theoretic bounds on the required interaction steps for adaptation could guide practical deployments and set appropriate expectations for performance in novel settings. Current personalization methods also implicitly assume alignment between local and global objectives, which may not hold in real-world applications where environment-specific rewards diverge significantly. Future algorithmic developments could incorporate negotiation mechanisms inspired by cooperative game theory to balance individual and collective goals in these multi-objective optimization scenarios.

Meta-learning frameworks represent another promising direction, where the federated model learns initialization parameters specifically designed to facilitate rapid personalization across diverse environments. Such approaches could be complemented by dynamic resource allocation strategies that prioritize training on environments where personalization yields the most significant marginal benefits to the global model. These advances would address the fundamental tension between personalization and generalization that characterizes federated reinforcement learning systems.

## 3.11 Scalability, Communication, and Implementation Challenges

The experiments documented in Jin et al. primarily focus on tabular and low-dimensional control tasks, leaving open questions about scaling FedRL to high-dimensional spaces such as pixel-based observations Jin et al. (2022b). This scaling introduces significant challenges related to the curse of dimensionality, as aggregating high-dimensional policies or value functions may require sophisticated compression techniques like knowledge distillation or network pruning to remain computationally feasible. The interaction between environment heterogeneity and partial observability represents another unexplored frontier, as

federated variants of POMDP algorithms have received minimal attention in the literature thus far. While DDPGAvg has been proposed for continuous action spaces, its efficacy in heterogeneous continuous control tasks such as multi-robot systems with varying dynamics requires more rigorous evaluation before deployment in complex domains.

Frequent model aggregation in FedRL introduces substantial communication costs, particularly when utilizing large neural networks. This challenge has prompted exploration of several potential solutions to improve communication efficiency. Adaptive communication schedules could allow agents to update their models locally until a performance plateau is detected, thereby reducing redundant transmissions without compromising learning progress. Delta parameter synchronization, which transmits only model parameter changes rather than full weights, has shown promise in federated supervised learning and could be adapted to the reinforcement learning context. Federated ensemble methods offer yet another approach, training smaller specialized models for each environment and aggregating them through ensemble techniques such as Bayesian committee machines.

## 3.12 Privacy, Security, and Ethical Considerations

While Jin et al. prevent explicit trajectory sharing to preserve privacy, sophisticated adversaries might still infer environment properties from the aggregated models Jin et al. (2022b). Strengthening privacy guarantees requires more robust mechanisms. Differential privacy (DP) approaches involving noise injection during aggregation or local updates could enhance protection, though this may exacerbate the suboptimality issues already present in FedRL systems. Quantifying the precise trade-off between differential privacy guarantees and utility in the FedRL context represents a critical area for future research. Secure multiparty computation (SMPC) offers an alternative approach, implementing cryptographic protocols to perform aggregation without revealing individual model updates, though these come with significant computational overhead that must be carefully managed.

The distinction between cross-device and cross-silo federated RL introduces additional complexities. Most current FedRL research assumes cross-device settings with numerous simple agents, but cross-silo scenarios-involving fewer but more capable agents, such as autonomous vehicles or industrial robots-present distinct challenges. Heterogeneous computing capabilities across these agents necessitate asynchronous aggregation protocols to synchronize training effectively. Non-IID environment distributions in cross-silo settings, which may exhibit systematic heterogeneity due to factors like regional differences, could benefit from clustered aggregation strategies that recognize and leverage these patterns.

Deploying FedRL policies from simulation to physical systems introduces the challenge of accounting for real-world heterogeneity not fully captured during training. Promising approaches to address this sim-to-real gap include domain randomization within federated training, explicitly sampling environment pa-

rameters from broader distributions to improve robustness. Additionally, online adaptation protocols could continuously incorporate real-world interaction data while preserving privacy constraints, allowing for ongoing refinement of models in deployment settings. In open federated systems with self-interested agents, participation may be limited if collaboration offers minimal local benefits. Economic frameworks could address this through mechanisms that price personalized model improvements, allowing agents to "purchase" better personalization through increased participation in the federated learning process.

The environmental impact of training large-scale FedRL systems is another important consideration, as these processes require significant energy consumption. Research into energy-efficient federated training through sparser communications, model sparsification, or specialized green computing techniques represents an imperative direction for responsible advancement of the field.

## 3.13   Synthesis and Future Directions

Federated reinforcement learning represents a powerful paradigm for collaborative learning across distributed agents with private experiences, but current approaches face significant limitations in theoretical foundations, personalization mechanisms, scalability, and privacy guarantees. Advancing beyond naive averaging techniques toward more sophisticated aggregation mechanisms that account for environment similarities could significantly reduce heterogeneity-induced suboptimality. Similarly, developing stronger theoretical understandings of embedding-based personalization approaches would enable more principled algorithm design and performance guarantees across diverse environments.

As the field moves toward high-dimensional observation spaces and more complex tasks, addressing communication efficiency and scalability becomes increasingly critical. Privacy and security considerations must evolve in parallel with algorithmic advances to ensure that the benefits of federated learning-particularly the privacy preservation of local trajectories-are not undermined by sophisticated inference attacks on aggregated models. Finally, ethical considerations around incentive design, and environmental impact should inform research directions to ensure that federated reinforcement learning technologiesmitigate benefit diverse stakeholders while minimizing potential harms.

By addressing these challenges through a coordinated research agenda, the federated reinforcement learning community can develop more robust, efficient, and equitable approaches to collaborative learning across heterogeneous environments, ultimately unlocking the full potential of this promising paradigm for applications ranging from robotics and autonomous vehicles to personalized healthcare and smart infrastructure systems.

# 4 Cooperative Decentralized RL

## 4.1 Introduction

Cooperative Decentralized Reinforcement Learning (CDRL) is an emerging research area that combines the principles of decentralized learning with the strategic framework of multi-agent reinforcement learning (MARL). In CDRL, multiple agents learn to act cooperatively without relying on a centralized controller, instead making decisions based on local observations and limited communication with nearby peers. This paradigm is particularly relevant for real-world applications where centralized control is impractical due to communication constraints, scalability issues, or privacy concerns. In decentralized learning, each agent maintains its own local policy and value estimates, yet must collaborate to achieve a shared objective. The inherent challenge lies in coordinating learning across distributed nodes with partial information, which often requires designing robust communication protocols and consensus mechanisms. MARL provides the formal framework for this setting by modeling the interactions among multiple agents, whose joint actions influence the environment and the resulting rewards. Cooperative MARL, a subset of this framework, specifically focuses on scenarios where all agents share a common goal, such as maximizing a global return. In these settings, the agents are incentivized to coordinate their actions, yet they must do so in a manner that preserves decentralized operation. The work by Zhang et al. (2018) represents a significant advancement in CDRL by addressing these challenges with a fully decentralized approach that leverages consensus updates and decentralized actor–critic algorithms.

In the following sections, we review the contributions of Zhang et al. (2018)'s paper in detail, examining how their methods effectively address the challenges of cooperative decentralized reinforcement learning.

## 4.2 Zhang et al (2018) Review

Zhang et al. (2018) propose a framework that enables agents to operate autonomously without the need for a centralized controller. In this setup, each agent is designed to independently process a local reward—which may differ from agent to agent—while determining its actions through an individualized, parameterized policy. Rather than relying on a global communication scheme, agents share information only with their immediate neighbors over a network whose connectivity may vary over time. Despite this variability, the system is structured such that the communication network remains jointly connected over bounded time intervals, ensuring that every agent can indirectly interact with all others.

The global objective is to maximize the *globally averaged return* defined as

$$J(\theta) = \sum_{s \in \mathcal{S}} d_\theta(s) \sum_{a \in \mathcal{A}} \pi_\theta(s, a) R(s, a),$$

where $d_\theta(s)$ is the stationary state distribution under the joint policy $\pi_\theta$ and

$$R(s, a) = \frac{1}{N} \sum_{i=1}^{N} R_i(s, a)$$

is the averaged reward over $N$ agents. Importantly, while the communication network is time-varying and need not be connected at every time step, it is assumed to be *jointly connected over time*—i.e., there exists a bounded time interval $T$ within which every agent can indirectly communicate with every other agent.
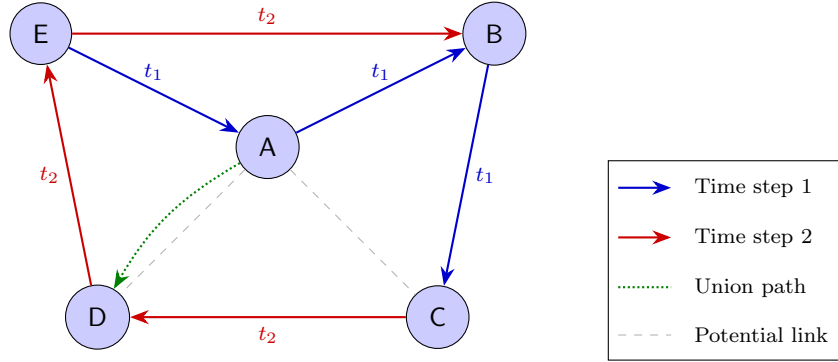


Figure 5: A dynamic communication topology illustrating time-varying connectivity among agents. Direct links are active at different time steps ($t_1$ and $t_2$), while dashed lines represent potential connections. The dotted green path shows an example of indirect connectivity achieved through the union of connections over time.

### 4.2.1 Single-Agent MDP and Actor–Critic Methods

A standard Markov decision process (MDP) is defined by the tuple $(\mathcal{S}, \mathcal{A}, P, r)$, where $\mathcal{S}$ represents a finite state space, $\mathcal{A}$ a finite action space, $P(s'|s, a)$ the state transition probability, and $r(s, a)$ the expected reward for taking action $a$ in state $s$. The objective in such a setting is to maximize the long-term average reward given by

$$J(\pi) = \sum_{s \in \mathcal{S}} d_\pi(s) \sum_{a \in \mathcal{A}} \pi(s, a) r(s, a),$$

where $d_\pi(s)$ denotes the stationary distribution under the policy $\pi$.

Actor–critic (AC) methods decompose the problem into:

(a) A *critic* that estimates the value (or action-value) function via function approximation.

(b) An *actor* that updates the policy parameters through gradient ascent.

The classical policy gradient theorem states that

$$\nabla_\theta J(\theta) = \mathbb{E}_{s \sim d_\theta,\, a \sim \pi_\theta} \left[ \nabla_\theta \log \pi_\theta(s, a)\, A_\theta(s, a) \right],$$

with the advantage function $A_\theta(s, a) = Q_\theta(s, a) - V_\theta(s)$.

### 4.2.2 Networked Multi-Agent MDP (NM-MDP)

Zhang et al. (2018) extend the single-agent Markov Decision Process (MDP) to a multi-agent framework involving $N$ agents, which they refer to as the Networked Multi-Agent MDP (NM-MDP). In this setting, each agent $i$ operates within its own local action space $A_i$ and adheres to a local policy $\pi_{\theta^i}^i(s, a^i)$. The overall joint policy is constructed by the factorization

$$\pi_\theta(s, a) = \prod_{i=1}^{N} \pi_{\theta^i}^i(s, a^i),$$

with the joint parameter vector defined as $\theta = [\theta^1, \ldots, \theta^N]$. Under this joint policy, the system's behavior is characterized by the stationary state distribution $d_\theta(s)$ and an averaged reward given by

$$R(s, a) = \frac{1}{N} \sum_{i=1}^{N} R_i(s, a).$$

Key assumptions in the NM-MDP framework include several conditions to ensure theoretical soundness. The communication graph $\mathcal{G}_t = (\mathcal{N}, \mathcal{E}_t)$ is allowed to vary over time, yet it is assumed to be jointly connected over time so that every agent can eventually interact with any other despite intermittent disconnections. Moreover, the induced Markov chain under any joint policy $\pi_\theta$ is assumed to be ergodic, ensuring long-term statistical regularity. The instantaneous rewards $r_t^i$ are uniformly bounded to avoid issues with unbounded variations. The consensus weight matrices $C_t = [c_t(i, j)]$ satisfy the condition

$$\sum_{j \in \mathcal{N}} c_t(i, j) = 1 \quad \forall\, i,$$

with $c_t(i, j) = 0$ whenever $(i, j) \notin \mathcal{E}_t$. In addition, in expectation, these matrices are column-stochastic, and the spectral norm of

$$E \left[ C_t^\top \left( I - \frac{1}{N} \mathbf{1} \mathbf{1}^\top \right) C_t \right]$$

is strictly less than one.

### 4.2.3 Decentralized Policy Gradient Theorem

Zhang et al. (2018) derive a decentralized policy gradient theorem using the log-derivative trick. Starting from the standard policy gradient,

$$\nabla_\theta J(\theta) = \mathbb{E}_{s, a} \left[ \nabla_\theta \log \pi_\theta(s, a)\, Q_\theta(s, a) \right],$$

and noting the factorization

$$\pi_\theta(s, a) = \prod_{i=1}^{N} \pi_{\theta^i}^i(s, a^i),$$

one obtains for each agent $i$:

$$\nabla_{\theta^i} J(\theta) = \mathbb{E}_{s,a} \left[ \nabla_{\theta^i} \log \pi_{\theta^i}^i(s, a^i) \, Q_\theta(s, a) \right].$$

By introducing the *local advantage function*

$$A_\theta^i(s, a) = Q_\theta(s, a) - \widetilde{V}_\theta^i(s, a^{-i}),$$

with

$$\widetilde{V}_\theta^i(s, a^{-i}) = \sum_{a^i \in A_i} \pi_{\theta^i}^i(s, a^i) \, Q_\theta(s, a^i, a^{-i}),$$

the gradient becomes

$$\nabla_{\theta^i} J(\theta) = \mathbb{E}_{s \sim d_\theta, \, a \sim \pi_\theta} \left[ \nabla_{\theta^i} \log \pi_{\theta^i}^i(s, a^i) \, A_\theta^i(s, a) \right].$$

In practice, each agent obtains a local estimate of $Q_\theta(s, a)$ via consensus updates in the critic, ensuring that all agents eventually agree on a common approximation.

### 4.2.4 Decentralized Actor–Critic Algorithms

Zhang et al. (2018) propose two actor–critic algorithms that use consensus-based critic updates and two-time-scale updates.

**Algorithm 1: Actor–Critic with Action-Value Function Approximation**

Each agent approximates the global action-value function using linear function approximation:

$$Q(s, a; \omega) = \omega^\top \phi(s, a),$$

where $\phi(s, a)$ is the feature vector. The critic parameters are updated locally and then fused via consensus with weight matrix $C_t$.

**Key Steps:**

- **Local Reward and Running Average:** Each agent $i$ maintains an estimate $\mu_i$ of its average reward.

- **Critic Update:** For agent $i$, the temporal-difference error is computed as
$$\delta_i^t = r_{t+1}^i - \mu_i^t + Q(s_{t+1}, a_{t+1}; \omega_i^t) - Q(s_t, a_t; \omega_i^t).$$
A local critic update is performed:
$$\widetilde{\omega}_i^t = \omega_i^t + \beta_{\omega,t} \, \delta_i^t \, \nabla_\omega Q(s_t, a_t; \omega_i^t).$$

- **Consensus Update:** Each agent updates its critic parameters via

$$\omega_i^{t+1} = \sum_{j \in \mathcal{N}} c_t(i,j)\, \widetilde{\omega}_j^t.$$

- **Actor Update:** The local advantage is calculated as

$$A_i^t = Q(s_t, a_t; \omega_i^t) - \sum_{a_i \in A_i} \pi_{\theta_i}^i(s_t, a_i)\, Q(s_t, a_i, a_{-i}; \omega_i^t),$$

and the policy update is performed using the score function

$$\psi_i^t = \nabla_{\theta_i} \log \pi_{\theta_i}^i(s_t, a_i^t).$$

The policy is then updated via

$$\theta_i^{t+1} = \Gamma_i\Big(\theta_i^t + \beta_{\theta,t}\, A_i^t\, \psi_i^t\Big),$$

where $\Gamma_i$ is the projection operator onto the compact set $\Theta_i$ (as per Assumption 4.1).

- **Policy Parameterization:** Although the paper permits general differentiable parameterizations, a common example is thesoftmax function:

$$\pi_{\theta_i}^i(s, a_i) = \frac{\exp(f_{\theta_i}(s, a_i))}{\sum_{a_i'} \exp(f_{\theta_i}(s, a_i'))},$$

where $f_{\theta_i}(s, a_i)$ is differentiable.

**Algorithm 6** Networked Actor–Critic with Action-Value Function Approximation

---

1: **Input:** Initial values $\mu_i^0$, $\omega_i^0$, $\widetilde{\omega}_i^0$, $\theta_i^0$ for all $i \in \mathcal{N}$; initial state $s_0$; stepsizes $\{\beta_{\omega,t}\}$ and $\{\beta_{\theta,t}\}$.

2: Each agent $i$ executes $a_i^0 \sim \pi_{\theta_i^0}^i(s_0, \cdot)$ and observes joint actions $a_0 = (a_1^0, \dots, a_N^0)$.

3: Initialize $t \leftarrow 0$.

4: **repeat**

5:     **for** each agent $i \in \mathcal{N}$ **do**

6:         Observe $s_{t+1}$ and reward $r_{t+1}^i$.

7:         Update running average:

$$\mu_i^{t+1} \leftarrow (1 - \beta_{\omega,t})\mu_i^t + \beta_{\omega,t}\, r_{t+1}^i.$$

8:         Select action: $a_i^{t+1} \sim \pi_{\theta_i^t}^i(s_{t+1}, \cdot)$.

9:     **end for**

10:     Observe joint actions $a_{t+1} = (a_1^{t+1}, \dots, a_N^{t+1})$.

11:     **for** each agent $i \in \mathcal{N}$ **do**

12:         Compute TD-error:

$$\delta_i^t \leftarrow r_{t+1}^i - \mu_i^t + Q(s_{t+1}, a_{t+1}; \omega_i^t) - Q(s_t, a_t; \omega_i^t).$$

13:         **Critic Update:**

$$\widetilde{\omega}_i^t \leftarrow \omega_i^t + \beta_{\omega,t}\, \delta_i^t\, \nabla_\omega Q(s_t, a_t; \omega_i^t).$$

14:         Compute local advantage:

$$A_i^t \leftarrow Q(s_t, a_t; \omega_i^t) - \sum_{a_i \in A_i} \pi_{\theta_i}^i(s_t, a_i)\, Q(s_t, a_i, a_{-i}; \omega_i^t).$$

15:         Compute score function:

$$\psi_i^t \leftarrow \nabla_{\theta_i} \log \pi_{\theta_i}^i(s_t, a_i^t).$$

16:         **Actor Update:**

$$\theta_i^{t+1} \leftarrow \Gamma_i\Big(\theta_i^t + \beta_{\theta,t}\, A_i^t\, \psi_i^t\Big).$$

17:         Send $\widetilde{\omega}_i^t$ to neighbors $\{j \in \mathcal{N} : (i,j) \in \mathcal{E}_t\}$.

18:     **end for**

19:     **for** each agent $i \in \mathcal{N}$ **do**

20:         **Consensus Step:**

$$\omega_i^{t+1} \leftarrow \sum_{j \in \mathcal{N}} c_t(i,j)\, \widetilde{\omega}_j^t.$$

21:     **end for**

22:     $t \leftarrow t + 1$.

23: **until** Convergence

---

**Pseudocode:**

**Algorithm 2: Actor–Critic with State-Value TD-Error**

In Algorithm 2, each agent employs a state-value function $V(s; v)$ and a global reward estimator $R(s, a; \lambda)$. Consensus updates are applied to the running average $\mu_i$, the reward estimator $\lambda_i$, and the critic parameter $v_i$.

**Key Steps:**

- **Running Average and Reward Estimation:**

$$\mu_{e,i}^t \leftarrow (1 - \beta_{v,t})\mu_i^t + \beta_{v,t}\, r_{t+1}^i,$$

$$\lambda_{e,i}^t \leftarrow \lambda_i^t + \beta_{v,t}\Big[r_{t+1}^i - R_t(\lambda_i^t)\Big]\nabla_\lambda R_t(\lambda_i^t).$$

- **Critic Update:** Compute the state-value TD-error:

$$\delta_i^t \leftarrow r_{t+1}^i - \mu_i^t + V(s_{t+1}; v_i^t) - V(s_t; v_i^t),$$

and update:
$$v_{e,i}^t \leftarrow v_i^t + \beta_{v,t}\, \delta_i^t \, \nabla_v V(s_t; v_i^t).$$

- **Actor Update:** Define the adjusted TD-error:

$$\delta_{e,i}^t \leftarrow R_t(\lambda_i^t) - \mu_i^t + V(s_{t+1}; v_i^t) - V(s_t; v_i^t),$$

then update:

$$\theta_i^{t+1} \leftarrow \Gamma_i\Big(\theta_i^t + \beta_{\theta,t}\, \delta_{e,i}^t\, \nabla_{\theta_i} \log \pi_{\theta_i}^i(s_t, a_i^t)\Big).$$

- **Consensus Updates:** Fuse the estimates via:

$$\mu_i^{t+1} \leftarrow \sum_{j\in\mathcal{N}} c_t(i,j)\,\mu_{e,j}^t, \quad \lambda_i^{t+1} \leftarrow \sum_{j\in\mathcal{N}} c_t(i,j)\,\lambda_{e,j}^t, \quad v_i^{t+1} \leftarrow \sum_{j\in\mathcal{N}} c_t(i,j)\,v_{e,j}^t.$$

**Algorithm 7** Networked Actor–Critic with State-Value TD-Error
___
1: **Input:** Initial values $\mu_i^0$, $\mu_{e,i}^0$, $v_i^0$, $v_{e,i}^0$, $\lambda_i^0$, $\lambda_{e,i}^0$, $\theta_i^0$ for all $i \in \mathcal{N}$; initial state $s_0$; stepsizes $\{\beta_{v,t}\}$ and $\{\beta_{\theta,t}\}$.
2: Each agent $i$ executes $a_i^0 \sim \pi_{\theta_i^0}^i(s_0, \cdot)$.
3: Initialize $t \leftarrow 0$.
4: **repeat**
5:     **for** each agent $i \in \mathcal{N}$ **do**
6:         Observe $s_{t+1}$ and reward $r_{t+1}^i$.
7:         Update running average:
$$\mu_{e,i}^t \leftarrow (1 - \beta_{v,t})\mu_i^t + \beta_{v,t}\, r_{t+1}^i.$$

8:         Update reward estimator:
$$\lambda_{e,i}^t \leftarrow \lambda_i^t + \beta_{v,t}\left[ r_{t+1}^i - R_t(\lambda_i^t)\right]\nabla_\lambda R_t(\lambda_i^t).$$

9:         Compute state-value TD-error:
$$\delta_i^t \leftarrow r_{t+1}^i - \mu_i^t + V(s_{t+1}; v_i^t) - V(s_t; v_i^t).$$

10:         **Critic Update:**
$$v_{e,i}^t \leftarrow v_i^t + \beta_{v,t}\, \delta_i^t\, \nabla_v V(s_t; v_i^t).$$

11:         Compute adjusted TD-error:
$$\delta_{e,i}^t \leftarrow R_t(\lambda_i^t) - \mu_i^t + V(s_{t+1}; v_i^t) - V(s_t; v_i^t).$$

12:         Compute score function:
$$\psi_i^t \leftarrow \nabla_{\theta_i} \log \pi_{\theta_i}^i(s_t, a_i^t).$$

13:         **Actor Update:**
$$\theta_i^{t+1} \leftarrow \Gamma_i\left(\theta_i^t + \beta_{\theta,t}\, \delta_{e,i}^t\, \psi_i^t\right).$$

14:         Send $\mu_{e,i}^t$, $\lambda_{e,i}^t$, $v_{e,i}^t$ to neighbors over $\mathcal{G}_t$.
15:     **end for**
16:     **for** each agent $i \in \mathcal{N}$ **do**
17:         **Consensus Step:**
$$\mu_i^{t+1} \leftarrow \sum_{j\in\mathcal{N}} c_t(i,j)\, \mu_{e,j}^t, \quad \lambda_i^{t+1} \leftarrow \sum_{j\in\mathcal{N}} c_t(i,j)\, \lambda_{e,j}^t, \quad v_i^{t+1} \leftarrow \sum_{j\in\mathcal{N}} c_t(i,j)\, v_{e,j}^t.$$

18:     **end for**
19:     $t \leftarrow t + 1$.
20: **until** Convergence
___

**Pseudocode:**

### 4.2.5   Convergence Analysis

Zhang et al. (2018) establish the convergence of their decentralized learning algorithm under a series of standard assumptions. First, each agent's parameter vector, $\theta_i$, is projected onto a compact set $\Theta_i$ via the operator $\Gamma_i$. This projection guarantees that the parameters remain bounded, thereby ensuring system stability. Additionally, the authors assume that the instantaneous rewards, $r_t^i$, are uniformly bounded, which prevents divergence due to unbounded reward fluctuations.

Furthermore, the consensus mechanism is governed by the matrices $C_t = [c_t(i,j)]$, which are constructed to be row-stochastic and, in expectation, column-stochastic. A critical condition imposed is that the spectral norm of

$$\mathbb{E}\left[C_t^\top \left(I - \frac{1}{N}\mathbf{1}\mathbf{1}^\top\right) C_t\right]$$

remains strictly less than one. This requirement is essential for the stability and proper functioning of the consensus process across the network.

The convergence analysis also relies on a two-time-scale update approach. In this framework, the critic's step sizes (e.g., $\beta_{\omega,t}$ or $\beta_{v,t}$) are set to dominate the actor's step sizes ($\beta_{\theta,t}$), ensuring that the critic converges more rapidly than the actor. Finally, to secure theoretical guarantees, the critic employs linear function approximation techniques—such as

$$Q(s,a;\omega) = \omega^\top \phi(s,a)$$

or

$$V(s;v) = v^\top \varphi(s),$$

with the additional assumption that the corresponding feature matrices possess full column rank.

**Critic Convergence (Theorem 1)**   For a fixed joint policy $\pi_\theta$, the consensus-based critic updates form a contraction mapping. In particular, the critic estimates $\{\omega_i^t\}$ converge almost surely to a common limit $\omega_\theta$ that satisfies the projected Bellman equation:

$$\Phi^\top D_\theta^{s,a}\left(T_\theta^Q(\Phi\omega_\theta) - \Phi\omega_\theta\right) = 0,$$

where $T_\theta^Q$ is the Bellman operator for the action-value function and $D_\theta^{s,a}$ is a diagonal matrix with entries $d_\theta(s)\pi_\theta(s,a)$.

**Actor Convergence (Theorem 2)**   Using two-time-scale stochastic approximation and ODE-based analysis (e.g., Borkar's Theorem), Zhang et al. (2018)

show that the actor updates converge almost surely to asymptotically stable equilibria of the ODE

$$\dot{\theta}_i = \hat{\Gamma}_i \Big( \mathbb{E}_{s \sim d_\theta, a \sim \pi_\theta} \left[ A_\theta^i(s, a) \nabla_{\theta_i} \log \pi_{\theta_i}^i(s, a^i) \right] \Big),$$

where $\hat{\Gamma}_i$ denotes the differential of the projection operator $\Gamma_i$.

### 4.2.6    Discussion and Final Remarks

The work of Zhang et al. (2018) makes a substantial contribution to decentralized multi-agent reinforcement learning by introducing a Networked Multi-Agent MDP (NM-MDP) that accommodates time-varying, yet jointly connected, communication networks. The authors derive a decentralized policy gradient theorem using the log-derivative trick and function factorization, which underpins their novel approach.

Building on this theoretical foundation, two actor–critic algorithms are proposed. In these algorithms, the critic updates are achieved through complete consensus-based mechanisms, while the actor updates are maintained via projection-based methods. This combination not only ensures stability but also facilitates efficient learning in a decentralized environment.

Furthermore, convergence guarantees are rigorously established through fixed-point analysis and two-time-scale ordinary differential equation (ODE) methods. Experimental evaluations validate that the proposed decentralized algorithms perform nearly on par with centralized methods in terms of cumulative return. Additional advantages of the approach include enhanced privacy—since agents do not share their raw rewards or policies—and improved communication efficiency.

Looking ahead, the authors outline promising avenues for future research, including extensions to competitive settings, the incorporation of continuous actions, and the exploration of more general (e.g., nonlinear) function approximators.

## 4.3    Subsequent works

### 4.3.1    State of the Art and Current Challenges in Decentralized Cooperative MARL

Decentralized cooperative Multi-Agent Reinforcement Learning (MARL) has evolved into a pivotal paradigm for coordinating intelligent agent collectives in complex, distributed environments. The current landscape is marked by notable achievements, yet significant challenges remain, especially concerning real-world applicability and scalability to large agent populations. A central challenge lies in **scalability and complexity**. As agent counts increase, the joint action space and state space expand exponentially, rendering traditional centralized approaches computationally intractable. Decentralized methodologies offer a necessary alternative, yet even these struggle to maintain learning

efficiency and coordination efficacy in massively multi-agent systems with complex interaction dynamics. Compounding this is the issue of **communication constraints**. Practical deployments invariably involve limitations on communication bandwidth, network reliability, and latency. Algorithms must be designed to be resilient to these constraints, capable of learning and coordinating effectively even with intermittent or delayed information exchange Zhao et al. (2020). **Partial observability and non-stationarity** present further intertwined difficulties. In decentralized settings, agents operate with localized perceptions, lacking complete knowledge of the global state and other agents' actions. This inherent partial observability, coupled with the non-stationarity arising from the simultaneous, interdependent learning processes of multiple agents, creates a highly complex and unstable learning environment. Addressing these intertwined challenges is paramount for advancing the field towards practical and robust decentralized cooperative MARL solutions.

Beyond these core issues, the field is increasingly grappling with the complexities of **directed and dynamic communication topologies**. Real-world agent networks are often characterized by directed communication graphs, where information flow is inherently asymmetric. Moreover, these communication structures are frequently not static; they evolve over time due to agent mobility, network fluctuations, or task requirements. Algorithms must be developed to explicitly account for and adapt to these dynamic directed communication patterns. A related, yet distinct, challenge lies in leveraging **asynchronous and event-triggered communication**. Synchronous communication protocols, while simplifying analysis, are often inefficient and impractical for large-scale distributed systems. Asynchronous and event-triggered mechanisms, where communication occurs only when necessary or at irregular intervals, offer the potential for significant efficiency gains and reduced communication overhead Nabli et al. (2023); Zhao et al. (2020). However, the design of MARL algorithms that can effectively harness these communication paradigms is still in its early stages. Finally, despite empirical progress, the **theoretical foundations and performance guarantees** for decentralized MARL algorithms remain underdeveloped. Establishing rigorous theoretical frameworks that can predict convergence behavior, guarantee stability, and characterize sample complexity is essential for building trust and enabling the widespread adoption of these techniques in critical applications.

### 4.3.2 Primary Advances in Decentralized Cooperative MARL

Recent years have witnessed notable strides in decentralized cooperative MARL, driven by the need to overcome the challenges outlined above. These advancements can be broadly categorized into several key areas, each contributing to the increasing sophistication and applicability of decentralized MARL techniques. One significant area of progress is the adoption of **Graph Neural Networks (GNNs) for communication and coordination**. GNNs provide a powerful framework for representing agent interactions and facilitating information exchange in decentralized systems. By leveraging the graph structure of agent

networks, GNNs enable agents to learn complex communication patterns, reason about network topology, and achieve more effective coordination.

Another crucial advancement is the development and refinement of **decentralized actor-critic algorithms**. Actor-critic methods are inherently well-suited to MARL due to their ability to handle non-stationarity and partial observability. Decentralized variants of actor-critic algorithms allow agents to learn policies and value functions based on local information and limited communication, providing a practical approach to learning in complex decentralized environments. Zhang et al Zhang et al. (2018) contributes to this area by proposing a distributed actor-critic algorithm specifically designed for MARL over directed graphs, addressing both decentralization and directed communication challenges. Furthermore, there is growing recognition of the importance of explicitly **handling directed communication graphs** in decentralized MARL. Moving beyond simplified undirected graph assumptions, researchers are developing algorithms that can operate effectively in scenarios with asymmetric information flow. These methods often involve adaptations of GNN-based communication or decentralized actor-critic frameworks to account for the directed nature of communication.

To enhance the efficiency and scalability of decentralized learning, **asynchronous decentralized optimization** techniques are increasingly being explored within MARL. Asynchronous algorithms allow agents to update their learning parameters based on locally available information, without requiring strict synchronization with all other agents. This reduces communication bottlenecks and can significantly accelerate training in large-scale systems Nabli et al. (2023). Lian et al. (2018) directly addresses this aspect, providing methodologies for asynchronous updates in decentralized optimization, which can be adapted and extended to the MARL context. Finally, the field is witnessing a growing number of **applications in networked systems**, ranging from wireless ad hoc networks to traffic management and swarm robotics. These real-world applications serve as both a driving force and a validation ground for decentralized MARL research. They highlight the practical relevance of addressing communication constraints, dynamic environments, and the need for robust and efficient algorithms Xu et al. (2024).

### 4.3.3 Directed Graphs in Decentralized MARL

The explicit consideration of directed graphs to model communication topologies represents a critical refinement in decentralized MARL research. Unlike simpler undirected graph models that assume symmetric information exchange, directed graphs more accurately capture the realities of many multi-agent systems where communication is inherently asymmetric or unidirectional, reflecting scenarios where influence and information flow are not necessarily reciprocal.

**Mathematical Representation and Recent Advances**

A directed graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ used in a multi-agent system consists of a set of agents $\mathcal{V} = \{1, 2, \ldots, N\}$ and a set of directed edges $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$. An edge $(i, j) \in \mathcal{E}$ signifies that agent $j$ can receive information from agent $i$, while the reverse is not necessarily true, thereby representing an asymmetric flow of information. This fundamental concept has spurred several notable advancements in the field. Algorithms are emerging that explicitly **learn asymmetric policies**, recognizing that agents in directed graphs may play different roles based on their position in the communication structure. For instance, agents with high out-degree (transmitting information to many) might adopt more exploratory policies, while agents with high in-degree (receiving information from many) might focus on exploitation based on aggregated information Wang et al. (2020).

Another significant development is seen in **influence-based methods**. By quantifying the influence of each agent—using centrality measures from graph theory such as PageRank or eigenvector centrality—the learning process can prioritize information from influential agents within the network Tang et al. (2020). This approach allows for a more nuanced integration of information, where the contributions of influential agents are given greater weight during decision-making. For acyclic directed graphs, **topological sorting techniques** provides a linear ordering of agents consistent with the direction of edges. This ordering can be exploited to design more structured communication and learning algorithms. For example, agents can be updated in topological order, ensuring that information from upstream agents is incorporated before updating downstream agents.

Finally, **adaptations of graph convolutional networks (GCNs)** originally designed for undirected graphs, are being explored to handle directed graphs Li et al. (2021). Modifications include using different weight matrices for incoming and outgoing edges, or employing attention mechanisms to selectively aggregate information from neighbors based on edge direction and importance. These adaptations allow agents to effectively process information from their directed neighbors and learn in complex directed communication environments.Collectively, these advancements underscore the importance of incorporating directed graph models into decentralized MARL. By reflecting the inherent asymmetries in real-world communication networks, these models pave the way for more sophisticated and effective multi-agent learning strategies.

**Open Problems and Research Directions for Directed Graph MARL**

Despite these advancements, challenges persist. **Scalability to large directed networks** remains a significant hurdle. As the number of agents and the complexity of the directed graph increase, efficient information aggregation and policy learning become computationally demanding. Furthermore, establishing **optimality guarantees** in directed graph MARL is inherently more difficult than in simpler settings. The asymmetry of information flow complicates

convergence analysis and makes it challenging to characterize the conditions under which decentralized algorithms can achieve globally optimal or even approximately optimal solutions Dibangoye et al. (2018). Open problems include developing more scalable GCN architectures for directed graphs, designing algorithms with provable convergence guarantees in directed settings, and exploring the impact of different directed graph topologies on learning performance and coordination.

### 4.3.4 Asynchronous Communication in Decentralized MARL

The shift from synchronous to asynchronous communication models represents a paradigm shift in decentralized MARL, moving away from idealized assumptions towards more realistic and scalable system designs. In synchronous systems, agents operate in lockstep, requiring global synchronization at each step of interaction and learning. This model is often impractical for large-scale, distributed systems where communication delays and failures are common. Asynchronous communication models, in contrast, allow agents to operate on *independent timescales*, making decisions and updating their policies based on locally available information and communication arrivals, without waiting for global synchronization Nabli et al. (2023).

**Core Principles and Key Developments**

The core principle of asynchronous systems is *decoupling agent operations from global time*. Agents proceed at their own pace, driven by local events and information availability.This fundamental shift has led to several key developments in asynchronous MARL:

One major advancement is the development of **delay-tolerant algorithms**. These algorithms are inherently more tolerant to communication delays and network latency Nguyen et al. (2020).Techniques such as delayed gradient updates, importance weighting, and robust aggregation methods are employed to mitigate the negative impact of outdated information in asynchronous updates. Another development involves adapting **experience replay mechanisms** for asynchronous settings. In such environments, experiences gathered by different agents can become temporally misaligned due to variations in update frequencies and communication delays. To address this, experience replay methods have been enhanced with strategies like temporal alignment and importance sampling, ensuring that the replay buffers maintain coherent and useful experiences for effective learning Munos et al. (2016). Additionally, the **Asynchronous Advantage Actor-Critic (A3C) algorithm**, originally conceived for single-agent reinforcement learning, has been extended to the multi-agent context Mnih et al. (2016). These decentralized A3C variants allow agents to asynchronously update their actor and critic networks, leveraging parallel computation and reducing synchronization bottlenecks. Another notable innovation is the move towards **clock-free coordination**. Some asynchronous MARL approaches aim for clock-free coordination, eliminating the need for any

global clock or synchronization mechanism. These methods rely on purely local interactions and asynchronous message passing to achieve coordination, further enhancing scalability and robustness. Finally, **gossip-based learning approaches** have been explored as a means of asynchronous information aggregation and policy dissemination. Inspired by distributed consensus protocols, gossip algorithms facilitate the iterative exchange of information among neighboring agents Mateos et al. (2019). This process gradually leads to the convergence of policies or value functions across the network. Together, these developments highlight the dynamic and resilient nature of asynchronous MARL systems, paving the way for more adaptive and robust multi-agent learning strategies.

### Open Problems and Research Directions for Asynchronous MARL

Despite the advantages of asynchronous communication, significant research challenges remain. Establishing **convergence guarantees** for asynchronous MARL algorithms is a major open problem. The lack of synchronization and the presence of delays complicate the theoretical analysis, and new analytical frameworks are needed to rigorously prove convergence and characterize performance Dibangoye et al. (2018). Furthermore, achieving **optimal coordination** in asynchronous systems is inherently more challenging than in synchronous settings. The decentralized and asynchronous nature of updates can lead to coordination failures or suboptimal joint policies. Open research questions include designing novel asynchronous algorithms with provable convergence and coordination guarantees, developing adaptive synchronization strategies that dynamically adjust the level of asynchronicity based on system conditions, and exploring the trade-off between asynchronicity and coordination performance in different MARL tasks and environments. The design of **adaptive event-trigger design** in asynchronous settings is also an open problem.

### 4.3.5 Event-Triggered Processes in Decentralized MARL

Event-triggered communication offers a fundamentally different approach to communication management in decentralized MARL, moving beyond periodic or continuous information exchange to selective communication triggered by specific events or conditions Zhao et al. (2020); Nabli et al. (2023). In contrast to time-triggered approaches where communication occurs at fixed intervals, event-triggered mechanisms aim to minimize communication overhead by transmitting information only when it is deemed necessary for maintaining performance or achieving coordination.

### Fundamental Mechanism and Recent Innovations

The fundamental mechanism of event-triggered approaches is based on defining *triggering conditions* that determine when communication should occur. These conditions are typically based on monitoring local agent states, policy changes,

or performance metrics. Communication is triggered only when these conditions are violated, indicating a need for information exchange with neighbors. Mathematically, a triggering condition for agent $i$ to communicate with agent $j$ at time $t$ can be represented as:

$$f(s_i(t), \pi_i(t), \mathcal{I}_i(t)) > \theta \tag{23}$$

where $s_i(t)$ is the local state of agent $i$, $\pi_i(t)$ is its policy, $\mathcal{I}_i(t)$ is the information agent $i$ has about its neighbors, $f$ is a triggering function, and $\theta$ is a threshold. Communication is triggered when the triggering function $f$ exceeds the threshold $\theta$, indicating a significant change or deviation that necessitates information exchange.

Recent innovations in event-triggered MARL include several noteworthy developments. One advancement is the emergence of **self-triggered approaches**, where each agent independently determines when to communicate based on its local information, without requiring coordination with other agents for triggering decisions Dimos and Johansson (2014). This further enhances decentralization and reduces communication overhead. Another innovation involves the application of **Lyapunov-based trigger design**. Lyapunov stability theory is being used to design event-triggering conditions that guarantee system stability and convergence while minimizing communication Nowzari and Cortés (2019). Lyapunov functions are used to monitor system performance and trigger communication only when necessary to maintain stability or improve performance. Additionally, **distributed event detection methods** where distributed algorithms are being developed , allow agents to collaboratively detect events that require communication Seyboth et al. (2013). These methods allow agents to collectively monitor system-wide conditions and trigger communication in a decentralized manner, enabling more sophisticated event detection and triggering mechanisms. In privacy-sensitive contexts, **privacy-preserving triggers** are implemented to safeguard sensitive information by minimizing the data exchanged during communication Gupta et al. (2021).Triggering conditions can be based on local information only, or employ privacy-preserving communication protocols to protect sensitive data during information exchange. Finally, **adaptive threshold mechanisms** have been introduced to dynamically adjust the triggering threshold $\theta$ based on system conditions or learning progress Li et al. (2023). For instance,the threshold can be increased when communication is costly or when learning is progressing well with less frequent communication, and decreased when coordination is critical or when learning is stagnating.

These innovations underscore the evolving nature of event-triggered communication in MARL, highlighting a concerted effort to balance efficient information exchange with system stability and privacy.

## Open Problems and Research Directions for Event-Triggered MARL

Despite the promise of event-triggered MARL, several open problems and research challenges need to be addressed. A key challenge is the **open problem of**

**adaptive event-trigger design**. Designing adaptive mechanisms that can dynamically optimize triggering thresholds and functions in response to changing environments and learning dynamics is a complex but crucial area for future research. Furthermore, establishing **optimality guarantees** for event-triggered MARL algorithms remains a significant challenge Dibangoye et al. (2018). The intermittent and data-driven nature of communication complicates the theoretical analysis, and new tools are needed to analyze convergence, stability, and performance bounds. **Security considerations** in event-triggered MARL are also important. Event-triggered communication patterns could potentially be exploited by adversaries to disrupt communication or extract sensitive information. Research is needed to investigate the security vulnerabilities of event-triggered MARL systems and develop robust and secure triggering mechanisms. Finally, **scalability to large networks** remains a concern, as even event-triggered communication mechanisms may incur significant overhead in very large-scale multi-agent systems. Developing highly scalable and communication-efficient event-triggered MARL algorithms is an ongoing research direction.

### 4.3.6 Integration of Directed Graphs, Asynchronous Communication, and Event-Triggered Processes

Recent research is increasingly exploring the integration of directed graphs, asynchronous communication, and event-triggered processes to create more sophisticated and efficient decentralized MARL algorithms. Combining these concepts can lead to synergistic benefits, enabling highly scalable, communication-efficient, and robust multi-agent systems. Integrating event-triggered communication with dynamic directed graphs is a particularly promising direction where, communication links are not only directed and dynamic but also event-triggered, further reducing communication overhead and adapting to changing network topologies and communication needs. Also, combining asynchronous updates with event-triggered communication in value decomposition-based MARL algorithms can lead to highly efficient and scalable decentralized learning Nabli et al. (2023). Agents can asynchronously update their local value functions and trigger communication only when significant changes occur in their local estimates or when coordination is required, minimizing communication while maintaining learning progress. A significant challenge in integrating these concepts is the theoretical analysis of convergence under combined conditions. Analyzing the convergence of MARL algorithms that simultaneously incorporate directed graphs, asynchronous communication, and event-triggered processes is complex, requiring new theoretical tools and frameworks. Establishing convergence guarantees and performance bounds for these integrated approaches is a crucial area for future theoretical research Dibangoye et al. (2018).

## 4.4 Research Challenges and Open Problems

In conclusion, while decentralized cooperative MARL has made significant strides, several overarching research challenges and open problems remain, particularly

concerning the integration of directed graphs, asynchronous communication, and event-triggered processes. **Scalability to large networks** is a persistent challenge across all three areas Zhang et al. (2018); Nabli et al. (2023); Zhao et al. (2020). Developing algorithms that can effectively handle massive multi-agent systems with complex communication topologies and limited communication resources is crucial for real-world deployment. **Difficulties in establishing optimality guarantees** also remain a significant hurdle Dibangoye et al. (2018). The decentralized nature of learning, coupled with communication constraints and complex network structures, makes it challenging to provide rigorous guarantees on the optimality of learned policies. The **open problem of adaptive event-trigger design** is particularly critical for realizing the full potential of event-triggered MARL. Designing robust and adaptive triggering mechanisms that can optimize communication efficiency without sacrificing learning performance is a key research direction. Finally, **security considerations** are becoming increasingly important as decentralized MARL systems are deployed in more complex and potentially adversarial environments. Addressing security vulnerabilities and designing robust and privacy-preserving decentralized MARL algorithms is a crucial area for future research. Addressing these challenges will pave the way for the next generation of decentralized cooperative MARL algorithms, enabling their widespread application in diverse domains.

# 5 Noncooperative MARL

## 5.1 Problem Statement

Let us recall the components of the standard RL formalism, however with a number of different agents each with their own private reward function:

1. **Environment** ($\mathcal{E}$): The environment provides feedback to the agent in the form of a reward signal and a new state after the agent takes an action. The environment encapsulates everything outside the agent that affects its decision-making process.

2. **Agent** ($i$): Each agent $i \in \{1, \ldots, N\}$ makes decisions by interacting with the environment. The objective of each agent is to learn a policy $\pi_i : \mathcal{S} \to \mathcal{A}_i$ that maximizes its long-term cumulative reward.

3. **State Space** ($\mathcal{S}$): The state space $\mathcal{S}$ represents all possible states of the environment. At time $t$, the environment is in state $s_t \in \mathcal{S}$.

4. **Action Space** ($\mathcal{A}_i$): Each agent $i$ has an individual action space $\mathcal{A}_i$, and the joint action of all agents is denoted as $\mathbf{a} = (a_1, a_2, \ldots, a_N) \in \mathcal{A} = \mathcal{A}_1 \times \mathcal{A}_2 \times \cdots \times \mathcal{A}_N$.

5. **Reward Function** ($R_i(s, \mathbf{a})$): Each agent receives a reward $R_i : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ based on the state and joint action. Unlike cooperative settings where all agents optimize a shared reward, in noncooperative MARL, each agent optimizes its own reward function.

6. **Transition Function** ($P(s'|s, \mathbf{a})$): The transition probability function defines how the state evolves based on the joint action:

$$P(s'|s, \mathbf{a}) = \Pr(s_{t+1} = s'|s_t = s, \mathbf{a}_t = \mathbf{a}).$$

7. **Discount Factor** ($\gamma$): The discount factor $\gamma \in [0, 1]$ determines the importance of future rewards. Higher values of $\gamma$ prioritize long-term rewards over immediate rewards.

In a single-agent reinforcement learning scenario, the agent interacts with the environment to maximize the expected sum of future rewards:

$$\pi^* = \arg\max_{\pi} \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t)\right],$$

where $R(s_t, a_t)$ is the immediate reward at time $t$, and $\pi^*$ is the optimal policy that maximizes this expectation.

In a multi-agent reinforcement learning (MARL) setting, each agent $i$ seeks to maximize its own cumulative reward:

$$\pi_i^* = \arg\max_{\pi_i} \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t R_i(s_t, \mathbf{a}_t)|\pi_i, \pi_{-i}\right],$$

where $\pi_{-i}$ represents the policies of all other agents. This dependence on other agents' strategies introduces a nonstationary environment, requiring specialized learning algorithms for stability and convergence.

## 5.2 Extending to Noncooperative Multi-Agent Reinforcement Learning

Noncooperative Multi-Agent Reinforcement Learning (MARL) extends the single-agent framework to scenarios involving multiple agents that act in the same environment.

Unlike cooperative MARL, where agents work together to optimize a shared goal, noncooperative MARL consists of self-interested agents with independent objectives. This setting introduces unique challenges:

1. **Reward Structure**: Each agent maximizes its own reward, which may lead to conflicts. In contrast to cooperative settings where agents optimize a shared utility, competition often results in adversarial interactions.

2. **Action-Value Dependencies**: The optimal action-value function $Q_i(s, a)$ for each agent depends on the policies of other agents, leading to a dynamic and nonstationary learning process.

3. **Exploration-Exploitation Dilemma**: Since other agents are also learning and adapting, an agent must explore effectively while exploiting learned policies to maximize rewards.

4. **Scalability**: As the number of agents increases, the dimensionality of the joint state-action space grows exponentially, making optimization computationally expensive.

5. **Decentralization**: Agents may have only partial observability, requiring decentralized learning techniques to make optimal decisions without global state access.

In a multi-agent setting, we have $N$ agents, each denoted as $A_1, A_2, \ldots, A_N$, with corresponding state spaces $\mathcal{S}_1, \mathcal{S}_2, \ldots, \mathcal{S}_N$, and action spaces $\mathcal{A}_1, \mathcal{A}_2, \ldots, \mathcal{A}_N$. Each agent's goal is to find a policy $\pi_i$ that maximizes its own expected cumulative reward, which is given by the sum of rewards over time:

$$\pi_i^* = \arg\max_{\pi_i} \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t \mathcal{R}_i(s_t, a_t)\right],$$

where $\mathcal{R}_i(s_t, a_t)$ is the reward received by agent $i$ when it takes action $a_t$ in state $s_t$.

The primary challenge in noncooperative MARL arises from the fact that the environment is dynamically altered by the actions of multiple agents, each of which is trying to maximize its own reward. This creates a **nonstationary** environment, where the reward and transition functions depend on the actions of all agents, not just the environment's response to the agent's actions.

## 5.3   Nash Equilibrium in Noncooperative MARL

In the context of noncooperative MARL, the notion of **Nash equilibrium** (NE) from game theory becomes highly relevant. A Nash equilibrium is a concept where, given the strategies of all other agents, no agent can unilaterally improve its own reward by changing its own strategy. In mathematical terms, for a set of strategies $\{\pi_1^*, \pi_2^*, \ldots, \pi_N^*\}$, a Nash equilibrium is defined by the following condition:

$$\pi_i^* = \arg\max_{\pi_i} \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t \mathcal{R}_i(s_t, a_t) \mid \pi_{-i}\right],$$

where $\pi_{-i}$ represents the policies of all other agents except agent $i$.

In a multi-agent system, each agent's policy must be a best response to the policies of the other agents. If all agents are playing their best responses simultaneously, the system reaches a Nash equilibrium, where no agent has an incentive to deviate from its strategy.

However, Nash equilibria in multi-agent environments are notoriously difficult to compute and can be non-unique or unstable. Moreover, in many cases, agents may converge to **suboptimal** equilibria, where the collective outcome is not optimal for any of the agents involved.

## 5.4 Variations of Nash Equilibrium in Non-Cooperative MARL

### 5.4.1 Generalized Nash Equilibrium

Traditional Nash Equilibrium (NE) assumes that each agent's strategy set is fixed and independent of others. However, in many *multi-agent reinforcement learning (MARL)* settings, agents operate under *shared constraints*, such as resource limitations, safety conditions, or regulatory requirements. This leads to the concept of *Generalized Nash Equilibrium (GNE)*, where an agent's feasible strategies depend on the actions of others Facchinei and Pang (2010); Başar and Olsder (1999).

**Definition of Generalized Nash Equilibrium** A *Generalized Nash Equilibrium (GNE)* is a set of policies $\pi^* = (\pi_1^*, \ldots, \pi_N^*)$ such that for every agent $i$,

$$V_i(\pi_i^*, \pi_{-i}^*) \geq V_i(\pi_i, \pi_{-i}^*), \quad \forall \pi_i \in \mathcal{F}_i(\pi_{-i}), \tag{24}$$

where $\mathcal{F}_i(\pi_{-i})$ represents the *feasible strategy set* of agent $i$, which depends on the policies of other agents. Unlike standard Nash equilibria, where all policies are independently chosen, in a GNE, each agent's strategy is constrained by *coupled feasibility constraints* that affect all agents.

**Examples in MARL**

1. **Resource-Constrained Optimization:** Consider *multi-agent network bandwidth allocation*, where each agent (user) selects a transmission rate $a_i$ while ensuring that the *total bandwidth does not exceed a limit*:

$$\sum_{i=1}^{N} a_i \leq B. \tag{25}$$

   Each agent's feasible actions depend on others, making this a *Generalized Nash Game* Rosen (1965).

2. **Multi-Agent Traffic Coordination:** Autonomous vehicles at an intersection must select speeds while *avoiding collisions*. The feasible strategies for one vehicle depend on the speeds chosen by others, requiring a GNE-based solution Parise et al. (2016).

3. **Competitive Market Equilibria:** In *multi-agent bidding* (e.g., electricity markets), agents place bids while ensuring that *total demand does not exceed supply*, leading to market-based GNE formulations Scutari et al. (2010).

**Solving GNE in MARL**  Unlike standard NE, computing GNE requires handling *joint feasibility constraints*. Common solution methods include:

1. **Projected Gradient Descent (PGD):** Ensures updates satisfy shared constraints by *projecting* gradient steps onto the feasible region Facchinei and Pang (2003a).

2. **Primal-Dual Methods:** Uses *Lagrangian multipliers* to enforce constraints dynamically Pang and Fukushima (2005).

3. **Variational Inequality (VI) Approaches:** Reformulates MARL equilibrium conditions as a *variational inequality problem*, ensuring a stable solution under joint constraints Facchinei and Pang (2003b).

### 5.4.2   Pure and Mixed Strategy Equilibria

In non-cooperative MARL, where strategic interaction is paramount, the distinction between deterministic and stochastic policies is critical:

- **Pure Strategy Nash Equilibrium**: Each agent employs a deterministic policy $\pi_i(s) = a_i \in \mathcal{A}_i$. In competitive environments, pure strategies may be predictable and exploitable, but computationally simpler to identify. A pure strategy equilibrium exists when:

$$V_i^{(\pi_i^*, \boldsymbol{\pi}_{-i}^*)}(s) \geq V_i^{(a_i', \boldsymbol{\pi}_{-i}^*)}(s) \quad \forall a_i' \in \mathcal{A}_i, \forall s \in \mathcal{S}, \forall i \tag{26}$$

- **Mixed Strategy Nash Equilibrium**: Agents randomize according to probability distributions $\pi_i(s) \in \Delta(\mathcal{A}_i)$. In non-cooperative settings, mixed strategies offer protection against exploitation through unpredictability. Nash's theorem guarantees existence in finite non-cooperative games.

For any state $s$, the Q-function in non-cooperative MARL captures the expected return when agent $i$ takes action $a_i$ while other agents follow $\boldsymbol{\pi}_{-i}$:

$$Q_i^{\boldsymbol{\pi}}(s, a_i, \boldsymbol{a}_{-i}) = R_i(s, a_i, \boldsymbol{a}_{-i}) + \gamma \sum_{s' \in \mathcal{S}} P(s' \mid s, a_i, \boldsymbol{a}_{-i}) V_i^{\boldsymbol{\pi}}(s') \tag{27}$$

### 5.4.3   $\varepsilon$-Nash Equilibrium

In practical non-cooperative MARL, computational limitations often necessitate approximations. An $\varepsilon$-Nash equilibrium represents a joint policy $\boldsymbol{\pi}^\varepsilon$ where no agent can unilaterally improve their value by more than $\varepsilon$:

$$V_i^{(\pi_i^\varepsilon, \boldsymbol{\pi}_{-i}^\varepsilon)}(s) \geq V_i^{(\pi_i', \boldsymbol{\pi}_{-i}^\varepsilon)}(s) - \varepsilon \quad \forall s \in \mathcal{S}, \forall \pi_i', \forall i \tag{28}$$

This concept is essential in non-cooperative MARL where agents may satisfice rather than optimize, particularly in complex, high-dimensional domains.

### 5.4.4 Local Nash Equilibrium

In continuous action spaces common in many non-cooperative MARL problems, the concept of local Nash equilibrium becomes relevant. A joint policy $\boldsymbol{\pi}^*$ is a local Nash equilibrium if there exists a neighborhood $\mathcal{N}(\boldsymbol{\pi}^*)$ such that:

$$V_i^{(\pi_i^*, \boldsymbol{\pi}_{-i}^*)}(s) \geq V_i^{(\pi_i', \boldsymbol{\pi}_{-i}^*)}(s) \quad \forall \pi_i' \in \mathcal{N}(\pi_i^*), \forall s \in \mathcal{S}, \forall i \tag{29}$$

This addresses the challenge of finding equilibria in continuous action spaces where global optimality guarantees may be infeasible.

### 5.4.5 Markov Perfect Equilibrium (MPE)

MPE extends Nash equilibrium to dynamic stochastic games with state transitions, which is the foundation of non-cooperative MARL. It requires that agents' strategies form a Nash equilibrium in every subgame:

$$\pi_i^*(s) \in \arg\max_{\pi_i} V_i^{(\pi_i, \boldsymbol{\pi}_{-i}^*)}(s) \quad \forall s \in \mathcal{S}, \forall i \tag{30}$$

For non-cooperative MARL, MPE captures the strategic reasoning where agents consider not only immediate rewards but also how their actions influence future state distributions and therefore future competitive interactions.

## 5.5 Mean-Field Games and Nash Equilibrium

### 5.5.1 Mean-Field Theory in Multi-Agent Systems

Mean-field theory is a mathematical framework that simplifies multi-agent interactions by replacing explicit pairwise interactions with an aggregate statistical representation. In large-scale multi-agent reinforcement learning (MARL), where direct modeling of each agent's interaction is computationally intractable, mean-field methods enable agents to reason about the average behavior of the population instead of tracking individual opponents Lasry and Lions (2007); Yang et al. (2018).

This approach is inspired by statistical physics, where interactions among particles in a system are approximated using macroscopic field equations rather than computing individual interactions Huang et al. (2006). The same principle applies in MARL, where agents approximate the collective effect of the entire population rather than considering each agent separately.

### 5.5.2 Definition of the Mean-Field Approximation

Given a population of agents $N$, a mean-field approximation assumes that an individual agent's dynamics are governed by an interaction with the mean distribution of other agents rather than specific individuals. The mean-field distribution $\mu_t$ captures the statistical representation of the system at time $t$:

$$\mu_t = \mathbb{E}[s_t^i, a_t^i] \tag{31}$$

where $s_t^i$ and $a_t^i$ denote the state and action of agent $i$ at time $t$, and the expectation is taken over the entire population Yang et al. (2018).

Instead of computing interactions with each agent explicitly, the transition dynamics of an individual agent are given by:

$$s_{t+1}^i = f(s_t^i, a_t^i, \mu_t, w_t) \tag{32}$$

where:

- $f(\cdot)$ represents the system dynamics.

- $w_t$ captures random perturbations (e.g., noise, environment stochasticity).

The mean-field distribution evolves according to the collective policies of all agents:

$$\mu_{t+1} = T(\mu_t, \pi_t) \tag{33}$$

where $T$ is a mean-field transition operator that updates $\mu_t$ based on the agents' policy $\pi_t$. This formulation enables agents to approximate the behavior of a large system without tracking every individual agent, making it computationally efficient for large-scale MARL Carmona and Delarue (2018).

### 5.5.3  Mean-Field Nash Equilibrium

A Mean-Field Nash Equilibrium (MFNE) occurs when each agent, given the mean-field distribution $\mu^*$, follows an optimal policy $\pi^*$ such that no unilateral deviation improves its expected return Lasry and Lions (2007).

The value function for an agent $i$, interacting with the mean-field distribution, is given by:

$$V_{\pi_i, \mu_i}(s_i) = \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t R_i(s_t^i, a_t^i, \mu_t) \Big| s_0^i = s_i\right] \tag{34}$$

where:

- $R_i(s_t^i, a_t^i, \mu_t)$ is the reward function for agent $i$, which depends on its own state $s_t^i$, action $a_t^i$, and the mean-field distribution $\mu_t$.

- $\gamma$ is the discount factor.

A Mean-Field Nash Equilibrium (MFNE) is achieved when:

$$V_{\pi_i^*, \mu_i^*}(s_i) \geq V_{\pi_i', \mu_i^*}(s_i), \quad \forall \pi_i', \forall s_i, \forall i \tag{35}$$

which means that given the equilibrium mean-field distribution $\mu^*$, no agent can improve its expected return by deviating from $\pi^*$.

Additionally, the mean-field distribution remains self-consistent under equilibrium policies:

$$\mu^* = \mathbb{E}[s_t^i, a_t^i | \pi_i^*] \tag{36}$$

This ensures that each agent's best response aligns with the collective behavior, making MFNE a powerful concept for modeling strategic interactions in large populations.

### 5.5.4 Applications in Large-Scale MARL

Mean-field approaches have been widely adopted in large-scale multi-agent systems, particularly in environments where direct agent-to-agent interactions are impractical due to computational constraints Yang et al. (2018).

**Competitive AI Benchmarks: Meta MMO** A recent benchmark, Meta MMO, serves as a structured testbed for mean-field learning in large-agent, non-cooperative MARL environments Choe et al. (2024). Unlike traditional small-scale multi-agent games, Meta MMO features over 100 agents competing in diverse mini-games, providing insights into how equilibrium strategies emerge in high-dimensional settings.

1. **Survival Mode**: Agents must balance foraging for resources while avoiding combat. Mean-field approximations help agents learn optimal survival strategies in a dynamically changing environment.

2. **Sandwich Mode**: Agents face both NPC threats and competing teams, requiring policies that focus on long-term survival rather than short-term adversarial tactics.

Meta MMO demonstrates that mean-field methods effectively scale to large agent populations while maintaining equilibrium stability Choe et al. (2024). The structured nature of the environment allows researchers to analyze how mean-field approximations impact emergent agent behaviors and strategic decision-making in competitive MARL.

### 5.5.5 Challenges and Open Problems

Despite its advantages, mean-field learning in MARL presents several challenges:

1. **Stationarity Assumption**: The mean-field approximation assumes that the distribution evolves smoothly, but in many real-world settings, abrupt strategic shifts can break this assumption.

2. **Convergence Guarantees**: While mean-field Nash equilibria exist under certain conditions, efficiently converging to them remains an active area of research.

3. **Heterogeneity**: Mean-field models often assume homogeneous agents, whereas real-world agents have diverse policies, goals, and constraints. Extending mean-field theory to heterogeneous populations is an ongoing challenge Carmona and Delarue (2018).

4. **Scalability to Real-World Systems**: While mean-field approximations work well in simulations, real-world applications (e.g., autonomous traffic control, financial markets) require more sophisticated mechanisms to handle uncertainty and adversarial behavior.

## 5.6    Cooperative vs. Noncooperative Settings

Multi-agent reinforcement learning (MARL) can be broadly categorized into **cooperative** and **noncooperative** settings, each presenting unique challenges and solution approaches.

In **cooperative MARL**, agents share a common objective and collaborate to optimize a joint reward function. Techniques such as **centralized training with decentralized execution (CTDE)** and **team-based learning** enable agents to learn shared strategies while executing decisions independently. Communication and information sharing among agents are common, facilitating global convergence to optimal strategies.

Conversely, in **noncooperative MARL**, agents are self-interested and seek to maximize their individual rewards, often at the expense of others. This competitive nature introduces conflicts and necessitates game-theoretic solutions such as **Nash equilibrium** and **Pareto optimality** to model agent interactions effectively. Unlike cooperative settings, where stability is often ensured through joint optimization, noncooperative settings pose significant challenges due to adversarial dynamics and strategic adaptation.

Key characteristics that distinguish noncooperative MARL include:

1. **Nonstationarity**: The environment becomes nonstationary as agents adapt their policies in response to others. Traditional RL techniques struggle in such settings as they assume a fixed transition model.

2. **Exploration and Exploitation**: Unlike single-agent RL, where exploration strategies focus solely on the environment, in multi-agent settings, exploration must also account for the presence of learning opponents. Agents must balance discovering new strategies while avoiding exploitation by adversaries.

3. **Policy Dependence**: The optimal policy of each agent depends on the evolving policies of others, creating an interdependent learning process that complicates convergence.

4. **Scalability**: As the number of agents increases, the joint state and action spaces grow exponentially, leading to the **curse of dimensionality**. Efficient learning mechanisms must mitigate this complexity.

5. **Convergence Issues**: Unlike cooperative settings, where agents converge towards a shared optimal solution, noncooperative MARL may result in unstable dynamics or cycling behaviors. Ensuring convergence to an equilibrium—if one exists—is an ongoing research challenge.

These distinctions underscore the fundamental differences between cooperative and noncooperative MARL. While cooperative MARL leverages shared rewards to drive collaboration, noncooperative MARL relies on strategic decision-making and competitive optimization, making learning dynamics more complex and unpredictable.

Noncooperative Multi-Agent Reinforcement Learning (MARL) extends traditional RL to environments with multiple interacting agents, each pursuing its own reward without considering the others' objectives. The challenge lies in the interdependence of agents' actions and the nonstationary dynamics they introduce into the environment. Nash equilibrium provides a foundational concept for understanding optimal strategies in such environments, but achieving equilibrium is often computationally intractable and unstable. Researchers continue to explore algorithms and techniques to address these challenges, such as decentralized learning, centralized training, and game-theoretic methods. As MARL applications continue to grow in complexity, from robotics to autonomous vehicles, addressing these foundational challenges will be essential for building reliable, scalable systems.

## 5.7 Existing Work and Algorithms in Noncooperative MARL

Noncooperative Multi-Agent Reinforcement Learning (MARL) presents unique challenges such as non-stationarity, equilibrium computation, and strategic adaptation. To address these, research has evolved into distinct algorithmic categories: value-based, policy-based, hybrid, and advanced methods.

Value-based methods, such as Minimax-Q Learning and Nash Q-Learning, extend Q-learning to multi-agent settings, focusing on optimal value functions under competitive interactions. However, they struggle with scalability and dynamic opponent behaviors. Policy-based methods, including MAPPO and COMA, optimize policies directly, improving convergence and stability in complex environments. Hybrid methods combine the strengths of value- and policy-based approaches, with techniques like Mean Field MARL approximating multi-agent interactions for better scalability. Advanced methods, such as Evolutionary Strategies and Meta-Learning, explore adaptability and generalization beyond traditional reinforcement learning frameworks.

This section categorizes and examines these methods, highlighting their contributions and mathematical foundations.

### 5.7.1 Value-Based Methods

**Markov Games (Stochastic Games) – Littman (1994a)** Michael L. Littman introduced *Markov Games*, also called *Stochastic Games*, as a generalization of Markov Decision Processes (MDPs) to multi-agent settings. These games model interactions where multiple agents' rewards depend on their joint actions. A Markov game for $N$ agents consists of the tuple:

$$(S, A_1, A_2, ..., A_N, P, R_1, R_2, ..., R_N) \tag{37}$$

where $S$ is the set of states, $A_i$ is the action space for agent $i$, $P : S \times A_1 \times A_2 \times ... \times A_N \to \Delta(S)$ defines the state transition probabilities, and $R_i : S \times A_1 \times A_2 \times ... \times A_N \to \mathbb{R}$ is the reward function for agent $i$. This formalization provides a structure to analyze both cooperative and competitive interactions in reinforcement learning.

**Minimax-Q Learning – Littman (1994b)** In adversarial environments, Littman proposed *Minimax-Q Learning*, which extends Q-learning to zero-sum Markov games. The core idea is to compute the minimax value of the Q-function:

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left[ r + \gamma \max_{\pi} \min_{a'} Q(s', a') - Q(s, a) \right] \tag{38}$$

where $\pi$ represents the policy of the learning agent, $a'$ is the action of the opponent, and $\gamma$ is the discount factor. The max-min operation ensures that the agent optimizes for the worst-case opponent response. This algorithm effectively learns optimal strategies in competitive environments by integrating game-theoretic principles into reinforcement learning.

**Non-Stationarity in Multi-Agent Learning – Bowling and Veloso (2002)** Bowling and Veloso examined the challenge of *non-stationarity*, where agents learn simultaneously, causing the environment to change dynamically. They analyzed conditions under which self-interested learners can converge to stable strategies. They highlighted that standard Q-learning updates may no longer be valid in these dynamic environments because:

$$Q_i(s, a) \neq \mathbb{E} \left[ r_i + \gamma \max_{a'} Q_i(s', a') \right] \tag{39}$$

since the policy of the opponent is not fixed but evolving, making the standard Bellman update unstable. This realization motivated research into equilibrium-based learning.

**Nash Q-Learning – Hu and Wellman (2003)** Hu and Wellman proposed *Nash Q-Learning*, which extends Q-learning to *general-sum games* by incorporating Nash equilibrium computation at each step. Instead of computing a single max operator, the algorithm updates Q-values using Nash equilibria in each state:

$$Q_i(s, a_1, ..., a_N) \leftarrow (1-\alpha)Q_i(s, a_1, ..., a_N) + \alpha \left[ r_i + \gamma \sum_{s'} P(s'|s, a_1, ..., a_N)V_i(s') \right] \tag{40}$$

where $V_i(s)$ is computed as the Nash equilibrium value of the game at state $s$:

$$V_i(s) = \text{NashEquilibrium}(Q_i(s, \cdot)) \tag{41}$$

This ensures that each agent considers the equilibrium response of others, allowing for better learning in competitive and cooperative settings.

### 5.7.2  Policy-Based Methods

**Multi-Agent Deep Deterministic Policy Gradient (MADDPG) – Lowe et al. (2017)**  Lowe et al. developed *MADDPG*, which integrates deep reinforcement learning with actor-critic methods. The key idea is to use *centralized training* while keeping execution decentralized. The policy update for each agent follows:

$$\nabla_{\theta_i} J(\theta_i) = \mathbb{E}_{s,a \sim D} \left[ \nabla_{\theta_i} \pi_i(a_i|s) \nabla_{a_i} Q_i(s, a_1, ..., a_N) \right] \tag{42}$$

where $\pi_i(a_i|s)$ is the policy for agent $i$, $Q_i(s, a_1, ..., a_N)$ is the centralized Q-function considering all agents' actions, and $D$ is the replay buffer storing past experiences. The advantage of MADDPG is that it allows each agent to learn an optimal policy while leveraging centralized training to address non-stationarity.

paragraphMulti-Agent Proximal Policy Optimization (MAPPO) – Yu et al. (2021)

MAPPO extends the widely used Proximal Policy Optimization (PPO) framework to multi-agent settings, incorporating centralized value function training while maintaining decentralized execution. The key optimization update follows:

$$\mathcal{L}^{CLIP}(\theta) = \mathbb{E}_t \left[ \min \left( r_t(\theta) A_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) A_t \right) \right] \tag{43}$$

where $r_t(\theta)$ is the probability ratio between the new and old policy, and $A_t$ is the advantage function. MAPPO has been shown to stabilize training and improve performance over independent PPO baselines in cooperative and competitive multi-agent settings.

**Counterfactual Multi-Agent Policy Gradients (COMA) – Foerster et al. (2018b)**  COMA introduces a centralized critic to evaluate the counterfactual advantage of each agent's action in order to address credit assignment in cooperative MARL. The key advantage function is computed as:

$$A^i(s, a^i) = Q(s, a) - \sum_{a'^i} \pi^i(a'^i|s) Q(s, (a^{-i}, a'^i)) \tag{44}$$

This counterfactual baseline helps the learning process by reducing variance and improving sample efficiency in cooperative environments.

**Multi-Agent Trust Region Policy Optimization (MA-TRPO) – Wen et al. (2021)**  MA-TRPO extends the Trust Region Policy Optimization (TRPO) method to multi-agent systems by enforcing policy updates within a trust region:

$$\max_{\theta} \mathbb{E}\left[\frac{\pi_\theta(a|s)}{\pi_{\theta_{\text{old}}}(a|s)}A(s,a)\right], \quad \text{s.t. } D_{KL}(\pi_\theta, \pi_{\theta_{\text{old}}}) \leq \delta \tag{45}$$

where $D_{KL}$ represents the KL-divergence between the old and new policies, and $\delta$ controls the size of the policy update. MA-TRPO improves stability in learning by constraining drastic policy updates.

### 5.7.3 Hybrid and Advanced Methods

Hybrid methods combine value-based and policy-based techniques, while advanced methods integrate game-theoretic principles, meta-learning, or curriculum learning.

**Soft Actor-Critic for Multi-Agent Learning (MASAC) – Iqbal and Sha (2021)** MASAC extends Soft Actor-Critic (SAC) to the multi-agent domain by introducing entropy regularization to improve exploration:

$$J(\pi) = \mathbb{E}\left[Q(s,a) - \alpha \log \pi(a|s)\right] \tag{46}$$

where $\alpha$ is the temperature coefficient controlling exploration. MASAC enhances sample efficiency and robustness in stochastic multi-agent environments.

**Mean Field Reinforcement Learning – Yang and Wang (2018)** Mean-Field RL approximates interactions in large multi-agent systems by considering the mean effect of other agents rather than individual agent actions:

$$Q^i(s, a^i, \bar{a}^{-i}) \leftarrow Q^i(s, a^i, \bar{a}^{-i}) + \alpha \left[r^i + \gamma \max_{a'} Q^i(s', a', \bar{a}^{-i}) - Q^i(s, a^i, \bar{a}^{-i})\right] \tag{47}$$

where $\bar{a}^{-i}$ represents the average action of other agents. This approach is particularly useful for large-scale multi-agent systems such as traffic management and swarm robotics.

**Meta-Learning in Multi-Agent Systems – Finn et al. (2017)** Meta-learning (learning to learn) has been applied in multi-agent reinforcement learning to enable fast adaptation to new tasks. One prominent approach is MAML (Model-Agnostic Meta-Learning), which optimizes policies across multiple tasks:

$$\theta^* = \theta - \alpha \nabla_\theta \sum_{\mathcal{T}i} \mathcal{L}(\pi\theta - \beta \nabla_\theta \mathcal{L}(\mathcal{T}_i)) \tag{48}$$

This formulation allows agents to generalize quickly to new environments by leveraging prior experience.

**Evolutionary Strategies for MARL – Pourchot and Sigaud (2018)**
Evolutionary algorithms such as Evolution Strategies (ES) have been applied to
MARL by optimizing policies using population-based search:

$$\theta_{t+1} = \theta_t + \alpha \sum_{i=1}^{N} w_i \nabla_\theta f(\theta + \sigma \epsilon_i) \tag{49}$$

where $\sigma$ is the noise scale and $w_i$ are weights assigned based on policy fitness.
ES-based approaches are particularly useful in high-dimensional, multi-agent
coordination problems.

## 5.8 Challenges and Future Directions

### 5.8.1 Challenges in Non-Cooperative MARL

Non-cooperative Multi-Agent Reinforcement Learning (MARL) presents sev-
eral fundamental challenges that hinder efficient learning and stable equilibrium
computation. These challenges arise due to the strategic interactions between
self-interested agents, leading to complex dynamics that traditional single-agent
reinforcement learning (RL) methods struggle to address.

1. **Equilibrium Computation Complexity** Computing Nash Equilibria
   (NE) in multi-agent systems is known to be **PPAD-complete** Daskalakis
   et al. (2009), making exact solutions computationally intractable for large-
   scale environments. Moreover, Nash equilibria may be non-unique or un-
   stable, leading to convergence issues where agents oscillate between sub-
   optimal strategies Shoham and Leyton-Brown (2007).

2. **Non-Stationarity and Adaptation** Non-cooperative MARL environ-
   ments are inherently non-stationary because each agent's policy evolves
   in response to others. This violates the Markov property assumed in tradi-
   tional RL, causing learned policies to become obsolete as opponents adjust
   their strategies Hernandez-Leal et al. (2019a).

3. **Exploration-Exploitation Trade-off** Agents must explore opponent
   strategies while avoiding exploitation. Excessive exploration can lead to
   performance degradation when facing strategic adversaries. Algorithms
   must balance discovering new strategies with ensuring robust decision-
   making against adversaries Foerster et al. (2018c).

4. **Curse of Dimensionality and Scalability** As the number of agents
   increases, the state and action spaces grow exponentially, leading to the
   **curse of dimensionality** Yang and Wang (2018). Traditional RL tech-
   niques struggle with this exponential growth, making it difficult to scale
   MARL solutions effectively.

5. **Convergence and Stability of Learning** Ensuring convergence to an
   equilibrium in non-cooperative MARL remains an open problem. Many

learning algorithms suffer from cycling behavior, divergence, or convergence to suboptimal equilibria Mazumdar et al. (2020b).

### 5.8.2 Future Directions in Non-Cooperative MARL

Despite these challenges, several promising directions are emerging to improve the efficiency, stability, and adaptability of non-cooperative MARL.

1. **Scalable Equilibrium Approximation Methods**
   Given the computational intractability of exact Nash Equilibria Daskalakis et al. (2009), researchers have focused on developing scalable equilibrium approximation techniques. Approaches such as Mean-Field MARL Yang and Wang (2018) approximate the interactions among many agents by modeling their aggregate effect as a distribution, significantly reducing computational complexity. Furthermore, meta-learning strategies are being explored to accelerate convergence in equilibrium approximation Pfau et al. (2023).

2. **Learning in Non-Stationary Environments**
   To address non-stationarity, algorithms incorporating opponent modeling and adaptive learning have been introduced. Opponent-aware MARL Hernandez-Leal et al. (2019b) allows agents to adjust strategies dynamically based on observed behaviors, improving robustness in competitive settings. Additionally, Meta-MARL Lanctot (2017) leverages experience from past interactions to quickly adapt to new opponents, facilitating transfer learning in non-stationary multi-agent dynamics.

3. **Efficient Exploration Strategies**
   Advanced exploration techniques such as intrinsic motivation Pathak et al. (2017) and opponent-aware curiosity-driven exploration Foerster et al. (2018a) help agents discover better strategies while avoiding excessive risk. Recently, reinforcement learning frameworks integrating uncertainty estimation Osband et al. (2016) have shown promise in balancing exploration and exploitation in strategic interactions.

4. **Addressing the Curse of Dimensionality**
   To mitigate scalability issues, recent work has focused on factorized value functions Son et al. (2019) and graph-based MARL Jiang et al. (2020), where agents selectively model interactions with relevant neighbors rather than all agents. Hierarchical MARL Kulkarni et al. (2016) also decomposes decision-making into multiple levels, reducing computational burden while maintaining strategic depth.

5. **Convergence and Stability in Learning**
   Stabilizing learning dynamics remains a key research challenge. Techniques such as policy regularization Mazumdar et al. (2020a) and equilibrium propagation Balduzzi et al. (2018) aim to prevent oscillatory behavior and ensure convergence to stable solutions. Furthermore, integrating

meta-learning with equilibrium selection strategies Gao et al. (2022) has been explored to guide agents toward optimal equilibria in multi-agent systems.

## 5.9 Noncooperative MARL - Summary

Non-cooperative Multi-Agent Reinforcement Learning (MARL) presents a dynamic and complex landscape characterized by strategic interactions between self-interested agents. Unlike cooperative settings, where agents work towards a shared goal, non-cooperative MARL introduces challenges such as equilibrium computation complexity, non-stationarity, and exploration-exploitation trade-offs. The inherent adversarial and competitive nature of such environments makes learning stable and optimal policies significantly more difficult.

Recent advances have introduced scalable equilibrium approximation methods, opponent-aware learning strategies, and efficient exploration techniques to address these challenges. Mean-field MARL Yang and Wang (2018) has been particularly promising in mitigating scalability issues, while meta-learning approaches Pfau et al. (2023); Lanctot (2017); Gao et al. (2022) provide mechanisms for rapid adaptation in non-stationary environments. Additionally, methods such as hierarchical MARL Kulkarni et al. (2016) and graph-based interactions Jiang et al. (2020) have improved the efficiency of decision-making in high-dimensional multi-agent systems.

However, ensuring convergence to stable and desirable equilibria remains an open problem. Techniques like policy regularization Mazumdar et al. (2020a) and equilibrium selection strategies Balduzzi et al. (2018) show potential in stabilizing learning dynamics, but further research is required to refine these methods for large-scale applications. Addressing these open challenges will be critical for advancing non-cooperative MARL in real-world applications such as autonomous driving, strategic game-playing, and competitive robotics.

Future research should focus on integrating meta-learning, equilibrium refinement, and adaptive exploration mechanisms to build more robust and generalizable MARL solutions. By doing so, the field can move towards achieving more efficient, stable, and scalable learning in competitive multi-agent environments.

# 6    Conclusion

This work presented a comprehensive survey of multi agent RL, focusing on explication of the main formalisms and research considerations over a breadth of literature review. Within MARL, the three distinct paradigms, Federated, Cooperative, and Noncooperative environments indicate a single system with coordinating computing nodes, a set of disconnected agents communicating in a gossip network topology cooperatively, and a set of agents interacting in a network topology noncooperatively, respectively. While highlighting the algorithmic commonalities across these domains, we also indicate their structural

distinctions and associated computational desiderata. We present the state of the art as far as the most recent advances as well as contemporary open problems of interest. We hope that this technical summary presents a useful guide for researchers and workers in the field as well as those working in applied domains interested in incorporating MARL who are looking for guidance as to the appropriate tool for their application.

# References

Balduzzi, D., Racanière, S., Martens, J., and Foerster, J. N. (2018). The mechanics of n-player differentiable games. *International Conference on Machine Learning (ICML)*.

Başar, T. and Olsder, G. J. (1999). Dynamic noncooperative game theory. *SIAM Series in Applied Mathematics*.

Borkar, V. (2008). *Stochastic approximation. A dynamical systems viewpoint.*

Borkar, V. (2024). Stochastic approximation with two time scales: The general case.

Bowling, M. and Veloso, M. (2002). Multiagent learning using a variable learning rate. *Artificial Intelligence*, 136(2):215–250.

Carmona, R. and Delarue, F. (2018). Probabilistic theory of mean field games. *Springer*.

Chen, Y., Gu, B., Chen, K., Huang, Y., Su, Q., and Huang, Z. (2021). Federated reinforcement learning with localized value functions. In *International Conference on Machine Learning*, pages 1666–1676. PMLR.

Choe, K. W., Sullivan, R., and Suárez, J. (2024). Massively multiagent minigames for training generalist agents. *arXiv preprint*, arXiv:2406.05071.

Daskalakis, C., Goldberg, P. W., and Papadimitriou, C. H. (2009). The complexity of computing a nash equilibrium. *SIAM Journal on Computing*, 39(1):195–259.

Dibangoye, J. S., Alami, R., and Simonin, O. (2018). Learning to act in decentralized partially observable MDPs. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1).

Dimos, N. and Johansson, K. H. (2014). Decentralized event-triggered control and optimization. *IFAC Proceedings Volumes*, 47(3):981–986.

Doan, T. T. (2021). Nonlinear two-time-scale stochastic approximation: Convergence and finite-time performance.

Doan, T. T. and Romberg, J. (2019). Finite-time performance of distributed two-time-scale stochastic approximation.

Facchinei, F. and Pang, J.-S. (2003a). Finite-dimensional variational inequalities and complementarity problems. *Springer*, 1:1–60.

Facchinei, F. and Pang, J.-S. (2003b). *Finite-Dimensional Variational Inequalities and Complementarity Problems*. Springer.

Facchinei, F. and Pang, J.-S. (2010). Generalized nash equilibrium problems. *Annals of Operations Research*, 186:173–210.

Finn, C., Abbeel, P., and Levine, S. (2017). Model-agnostic meta-learning for fast adaptation of deep networks. *International Conference on Machine Learning (ICML)*.

Foerster, J. N., Chen, R. Y., Al-Shedivat, M., Whiteson, S., Abbeel, P., and Mordatch, I. (2018a). Learning with opponent-learning awareness. *International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*.

Foerster, J. N., Farquhar, G., Afouras, T., Nardelli, N., and Whiteson, S. (2018b). Counterfactual multi-agent policy gradients. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32.

Foerster, J. N., Farquhar, G., Afouras, T., Torr, P. H. S., and Whiteson, S. (2018c). Counterfactual multi-agent policy gradients. In *Proceedings of the 32nd AAAI Conference on Artificial Intelligence (AAAI-18)*, pages 2974–2982.

Gao, J., Zhang, Y., Xu, Z., and Lin, Y. (2022). Meta-learning for equilibrium selection in multi-agent reinforcement learning. *Artificial Intelligence Journal*, 306:1–25.

Gavi, K. (2025). Integrating reinforcement learning and neural networks for autonomous decision-making in embedded systems. *unknown*.

Gupta, P. L., Rabbat, R., and Savvides, M. (2021). Federated reinforcement learning with privacy-preserving event-triggered communication. *IEEE Transactions on Neural Networks and Learning Systems*, 33(11):6913–6927.

Hernandez-Leal, P., Kartal, B., and Taylor, M. E. (2019a). A survey and critique of multiagent deep reinforcement learning. *Autonomous Agents and Multi-Agent Systems*, 33:750–797.

Hernandez-Leal, P., Kartal, B., and Taylor, M. E. (2019b). A survey and critique of multiagent deep reinforcement learning. *Autonomous Agents and Multi-Agent Systems*, 33(3):750–797.

Hu, J. and Wellman, M. P. (2003). Nash q-learning for general-sum stochastic games. *Journal of Machine Learning Research*, 4:1039–1069.

Huang, M., Malhame, R. P., and Caines, P. E. (2006). Large-population cost-coupled lqg problems with nonuniform agents: Individual-mass behavior and decentralized $\varepsilon$-nash equilibria. *IEEE Transactions on Automatic Control*, 52(9):1560–1571.

Huh, D. and Mohapatra, P. (2024). Multi-agent reinforcement learning: A comprehensive survey.

Iqbal, S. and Sha, F. (2021). Randomized entity-wise factorization for multi-agent reinforcement learning. *International Conference on Machine Learning (ICML)*.

Jiang, J., Dun, C., Lu, W., Li, X., and Yu, Y. (2020). Graph-based multi-agent reinforcement learning: A survey. *arXiv preprint arXiv:2006.07281*.

Jin, H., Peng, Y., Yang, W., Wang, S., and Zhang, Z. (2022a). Federated reinforcement learning with environment heterogeneity.

Jin, H., Peng, Y., Yang, W., Wang, S., and Zhang, Z. (2022b). Federated reinforcement learning with environment heterogeneity. In *International Conference on Artificial Intelligence and Statistics*, pages 18–37. PMLR.

Kulkarni, T. D., Narasimhan, K., Saeedi, A., and Tenenbaum, J. B. (2016). Hierarchical deep reinforcement learning: Integrating temporal abstraction and intrinsic motivation. *Advances in Neural Information Processing Systems (NeurIPS)*, 29.

Lanctot, M. e. a. (2017). A unified game-theoretic approach to multiagent reinforcement learning. *NeurIPS*.

Lasry, J. M. and Lions, P. L. (2007). Mean-field games. *Japanese Journal of Mathematics*, 2:229–260.

Li, R., Wang, J., and Gao, Z. (2023). Adaptive event-triggered communication for multi-agent reinforcement learning in resource-constrained networks. *IEEE Internet of Things Journal*, 10(12):10908–10920.

Li, Y., Sun, J., Gong, M., and Grosse, R. (2021). Directed acyclic graph neural networks. In *International Conference on Learning Representations*.

Lian, X., Zhang, W., Zhang, C., and Liu, J. (2018). Asynchronous decentralized parallel stochastic gradient descent.

Littman, M. L. (1994a). Markov games as a framework for multi-agent reinforcement learning. In *Proceedings of the Eleventh International Conference on Machine Learning (ICML)*, pages 157–163.

Littman, M. L. (1994b). Minimax-q: Learning in zero-sum stochastic games. In *Proceedings of the Eleventh International Conference on Machine Learning (ICML)*, pages 157–163.

Lowe, R., Wu, Y., Tamar, A., Harb, J., Abbeel, P., and Mordatch, I. (2017). Multi-agent actor-critic for mixed cooperative-competitive environments. *arXiv preprint arXiv:1706.02275*.

Lyu, L., Jiang, Q., Yu, H., Yang, Q., and Yu, P. S. (2021). Privacy-preserving federated reinforcement learning for mobile crowdsensing. *arXiv preprint arXiv:2108.11887*.

Mateos, G., Faria, J., and Boyd, S. (2019). Gossip algorithms for decentralized multi-agent reinforcement learning. *IEEE Transactions on Signal Processing*, 67(19):5015–5028.

Mazumdar, E., Jordan, M. I., and Ratliff, L. J. (2020a). On gradient-based learning in continuous games. *SIAM Journal on Mathematics of Data Science*.

Mazumdar, E., Jordan, M. I., and Ratliff, L. J. (2020b). On the convergence of gradient-based learning in continuous games. *arXiv preprint*.

McMahan, B., Moore, E., Ramage, D., Hampson, S., and y Arcas, B. A. (2017). Communication-efficient learning of deep networks from decentralized data. *Artificial Intelligence and Statistics*, pages 1273–1282.

Mnih, V., Badia, A. P., Mirza, M., Graves, A., Lillicrap, T., Harley, P., Silver, D., and Kavukcuoglu, K. (2016). Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pages 1928–1937. PMLR.

Munos, R., Steinke, T., and Harutyunyan, A. (2016). Safe and efficient off-policy reinforcement learning. In *Advances in neural information processing systems*, pages 2939–2947.

Nabli, A., Belilovsky, E., and Oyallon, E. (2023). Accelerating asynchronous communication in decentralized deep learning. *Advances in Neural Information Processing Systems*, 36:47451–47474.

Nadiger, C., Kumar, A., and Abdelhak, S. (2019). Federated reinforcement learning for fast personalization. In *2019 IEEE Second International Conference on Artificial Intelligence and Knowledge Engineering (AIKE)*, pages 123–127. IEEE.

Naeem, M., Rizvi, S., and Coronato, A. (2020). A gentle introduction to reinforcement learning and its application in different fields. *IEEE Access*, 8:209320–209344.

Nguyen, Q. V., Phung, D., and Le, T. (2020). Optimistic asynchronous decentralized deep reinforcement learning. *arXiv preprint arXiv:2006.07485*.

Nowzari, C. and Cortés, J. (2019). Event-triggered reinforcement learning for networked control systems. *IEEE Control Systems Letters*, 4(2):243–248.

Osband, I., Russo, D. J., and Van Roy, B. (2016). Deep exploration via bootstrapped dqn. *Advances in Neural Information Processing Systems (NeurIPS)*, 29.

Pang, J.-S. and Fukushima, M. (2005). Quasi-variational inequalities, generalized nash equilibria, and multi-leader-follower games. *Computational Management Science*, 2:21–56.

Parise, F., Grammatico, S., and Lygeros, J. (2016). Network aggregative games: Distributed convergence to generalized nash equilibria. *IEEE Transactions on Automatic Control*, 61(11):3391–3406.

Pathak, D., Agrawal, P., Efros, A. A., and Darrell, T. (2017). Curiosity-driven exploration by self-supervised prediction. *International Conference on Machine Learning (ICML)*.

Pfau, D., Vinyals, O., and de Freitas, N. (2023). Meta-learning for equilibrium approximation in multi-agent systems. *Journal of Artificial Intelligence Research*, 76:1–24.

Pourchot, P. and Sigaud, O. (2018). Continuous control with deep reinforcement learning and evolution strategies. *Advances in Neural Information Processing Systems (NeurIPS)*.

Qi, J., Zhou, Q., Lei, L., and Zheng, K. (2021). Federated reinforcement learning: Techniques, applications, and open challenges. *arXiv preprint arXiv:2108.11887*.

Rosen, J. B. (1965). Existence and uniqueness of equilibrium points for concave n-person games. *Econometrica*, 33(3):520–534.

Scutari, G., Palomar, D. P., Facchinei, F., and Pang, J.-S. (2010). Convex optimization, game theory, and variational inequality theory. *IEEE Signal Processing Magazine*, 27(3):35–49.

Seyboth, G., Dimarogonas, D. V., and Johansson, K. H. (2013). Event-triggered distributed optimization. *IEEE Transactions on Automatic Control*, 58(2):246–262.

Shoham, Y. and Leyton-Brown, K. (2007). *Multiagent Systems: Algorithmic, Game-Theoretic, and Logical Foundations*. Cambridge University Press.

Son, K., Kim, D., Kang, W., Hostallero, D. E., and Yi, Y. (2019). Qtran: Learning to factorize with transformation for cooperative multi-agent reinforcement learning. *International Conference on Machine Learning (ICML)*.

Sutton, R. S. and Barto, A. G. (2018). *Reinforcement Learning: An Introduction*. MIT Press.

Tang, Y., Lou, T., and Shi, Y. (2020). Communication topology aware multi-agent reinforcement learning. In *IEEE INFOCOM 2020-IEEE Conference on Computer Communications*, pages 937–945. IEEE.

Truex, S., Baracaldo, N., Anwar, A., Steinke, T., Ludwig, H., Zhang, R., and Zhou, Y. (2019). A hybrid approach to privacy-preserving federated learning. In *Proceedings of the 12th ACM Workshop on Artificial Intelligence and Security*, pages 1–11.

Wang, P., Lu, Z., and Gao, Y. (2020). Tackling asymmetry via regularization in multi-agent reinforcement learning. *arXiv preprint arXiv:2007.14435*.

Wei, J., Huang, X., Ma, X., Jin, X., and Ren, F. (2020). Federated off-policy learning under distributed data with personalized privacy preservation. *Sensors*, 20(17):4883.

Wen, C., Devlin, S., Hofmann, K., and Rowland, M. (2021). Multi-agent trust region policy optimization. *International Conference on Autonomous Agents and Multi-Agent Systems (AAMAS)*.

Xu, Z., Zhu, Q., Sun, J., and Liu, H. (2024). Communication-efficient decentralized multi-agent reinforcement learning for cooperative adaptive cruise control. *arXiv preprint arXiv:2403.14952*.

Yang, Y. and Wang, J. (2018). Mean field multi-agent reinforcement learning. *International Conference on Machine Learning (ICML)*.

Yang, Y., Wang, J., Luo, R., and Li, M. (2018). Mean field multi-agent reinforcement learning. *International Conference on Machine Learning (ICML)*.

Yu, C., Zhan, X., and Yang, Y. (2021). The surprising effectiveness of MAPPO: Multi-agent proximal policy optimization. *Advances in Neural Information Processing Systems (NeurIPS)*.

Zeng, S., Doan, T. T., and Romberg, J. (2024a). A two-time-scale stochastic optimization framework with applications in control and reinforcement learning.

Zeng, S., Doan, T. T., and Romberg, J. (2024b). A two-time-scale stochastic optimization framework with applications in control and reinforcement learning. *SIAM Journal on Optimization*, 34(1):946–976.

Zhang, K., Yang, Z., Liu, H., Zhang, T., and Basar, T. (2018). Fully decentralized multi-agent reinforcement learning with networked agents. In Dy, J. and Krause, A., editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 5872–5881. PMLR.

Zhang, X., Chen, K., Zhang, W., and Xu, J. (2023). Advances in federated reinforcement learning: Methods and applications. *Applied Sciences*, 13(6):6497.

Zhao, M., Xu, W., Duan, J., and Lu, J. (2020). Event-triggered communication network with limited-bandwidth constraint for multi-agent reinforcement learning. *IEEE Transactions on Neural Networks and Learning Systems*, 32(12):5413–5427.

Zhuo, H. H., Feng, W., Lin, Y., Xu, Q., and Yang, Q. (2019). Federated deep reinforcement learning. *arXiv preprint arXiv:1901.08277*.