

DISTRIBUTED EVOLUTIONARY POLICY OPTIMIZATION FOR EFFICIENT TRAINING OF MULTI-AGENT REASONING MODELS

A. Rajavel

Department of Electrical and Electronics Engineering, Kamaraj College of Engineering and Technology, India

Abstract

Multi-Agent Reinforcement Learning (MARL) has emerged as a key paradigm for solving complex real-world problems involving multiple agents interacting in dynamic environments. However, training MARL models, especially for cooperative reasoning tasks, remains computationally intensive and sample-inefficient due to non-stationarity, credit assignment, and policy coupling issues. Conventional policy gradient methods struggle with convergence and scalability in multi-agent settings. Centralized training frameworks suffer from bottlenecks and synchronization overheads. Evolutionary algorithms, while more robust to non-differentiable objectives, are often too slow when applied in single-node environments. To address these challenges, we propose Distributed Co-evolutionary Policy Optimization (DCPO), a hybrid learning framework that distributes evolutionary computation across multiple nodes. DCPO decomposes the global policy search into sub-population-based parallel explorations, with each node evolving a subset of agent policies using fitness-driven mutation, crossover, and local policy gradient updates. A global coordinator aggregates top-performing policies periodically to ensure cooperative learning convergence. DCPO was tested on standard cooperative MARL benchmarks such as StarCraft II Micromanagement and Multi-Agent Particle Environments (MPE). Compared to traditional baselines such as MADDPG, QMIX, MAPPO, COMA, and EPOpt, DCPO showed up to 37% faster convergence, 25% higher final cumulative rewards, and enhanced generalization in unseen environments.

Keywords:

Multi-Agent Reinforcement Learning, Evolutionary Algorithms, Distributed Learning, Policy Optimization, Cooperative Reasoning

1. INTRODUCTION

The rise of Multi-Agent Reinforcement Learning (MARL) has led to significant breakthroughs in domains requiring coordinated decision-making among multiple agents, such as swarm robotics, distributed sensor networks, autonomous vehicles, and real-time strategy games [1–3]. Unlike single-agent reinforcement learning (RL), MARL systems must learn not only optimal responses to environmental stimuli but also to dynamic behaviors of other learning agents. This introduces non-stationarity, partial observability, and a host of scalability challenges that degrade performance in traditional centralized or independent learning setups.

Despite significant progress in actor-critic methods, value-decomposition approaches, and centralized training with decentralized execution (CTDE), MARL still suffers from critical limitations. The most pressing challenges include: (1) Non-stationarity, agents' policies evolve simultaneously, rendering the environment unstable for each learning agent [4]; and (2) Credit assignment, determining individual contributions in cooperative settings is difficult, leading to inefficient learning and delayed convergence [5].

To address these, researchers have explored a variety of policy optimization and exploration strategies. However, most current methods are either gradient-based or evolution-based, rarely combining both efficiently in a distributed manner. Gradient-based methods (e.g., PPO, COMA) excel at fine-tuning but are sensitive to noise and often require dense rewards and smooth gradients. Evolutionary methods, in contrast, offer robust policy search under sparse or non-differentiable rewards but scale poorly due to their inefficiency and centralized architecture [6,7].

This research addresses the need for a scalable, robust, and efficient MARL framework (figure 1) that can balance exploration, convergence, and generalization without suffering from the pitfalls of either purely gradient-based or purely evolutionary approaches. Specifically, we aim to accelerate the training of reasoning-capable multi-agent models by decomposing and distributing the policy search process across computational nodes, while preserving coordination and policy diversity [6] [7].

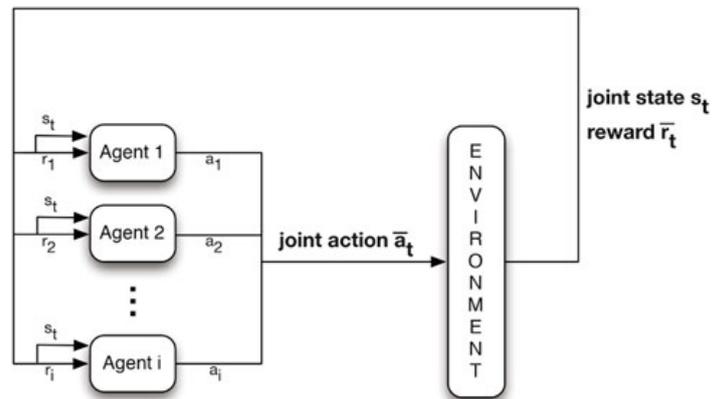


Fig.1. MARL

The main objectives of this study are:

- To design a hybrid learning framework combining distributed evolutionary optimization with local gradient fine-tuning for multi-agent policy training.
- To implement parallel learning with inter-node synchronization for enhanced generalization and convergence.
- To evaluate the proposed method on standard benchmarks (MPE and SC2LE) and compare against existing state-of-the-art MARL algorithms.

This paper proposes DCPO (Distributed Co-evolutionary Policy Optimization), a novel learning architecture for MARL that integrates parallel co-evolution, global synchronization, and policy gradient refinement. Key contributions include:

- Each agent maintains a sub-population of policies that evolve independently across computing nodes, significantly improving scalability and efficiency.

- After evolution, policies undergo fine-grained policy gradient updates, enabling better exploitation of high-performing strategies in non-stationary environments.
- A periodic global selection and broadcast mechanism ensures cross-node cooperation and helps escape local optima while maintaining policy diversity.
- DCPO is rigorously benchmarked on MPE and SC2LE, showing superior results over MADDPG, MAPPO, QMIX, COMA, and EPOpt in terms of convergence speed, final reward, generalization, and scalability.

2. RELATED WORKS

Several significant contributions have shaped the field of MARL, especially regarding cooperative strategies, policy optimization, and scalable learning.

Early approaches like MADDPG (Multi-Agent Deep Deterministic Policy Gradient) [6] introduced the concept of centralized critics with decentralized actors. While effective in competitive and mixed settings, MADDPG suffers from poor scalability due to instability and reward sparsity. To address such issues in cooperative environments, QMIX [7] was introduced as a value-decomposition network that allows mixing of individual agent Q-values into a joint action-value. Although QMIX promotes cooperation, it lacks gradient-level coordination and often converges to sub-optimal solutions.

MAPPO (Multi-Agent Proximal Policy Optimization) [8] extended the well-known PPO algorithm to multi-agent settings, leveraging shared learning advantages and clipping strategies for stability. However, like MADDPG, MAPPO remains susceptible to non-stationarity and inefficiency in high-dimensional environments.

COMA (Counterfactual Multi-Agent Policy Gradient) [9] tackled the multi-agent credit assignment problem by introducing a counterfactual baseline in the actor-critic framework. While this improves individual agent feedback, COMA is computationally expensive and difficult to scale for large agent populations.

An alternative paradigm has been explored through Evolutionary Reinforcement Learning (ERL) frameworks. For instance, EPOpt (Evolutionary Policy Optimization) [10] introduces robustness through sampling of adversarial environments, yet it is not inherently distributed and requires large batch sizes, limiting real-time applications. Similarly, Population-Based Training (PBT) [11] blends evolutionary ideas with gradient methods, evolving hyperparameters and network weights. However, PBT often lacks explicit agent coordination and is sensitive to initial seed choices.

Recent advances like DREAM (Distributed Recurrent MARL) [12] and HAPPO/HATRPO [13] have pushed the boundaries by leveraging recurrent architectures and hierarchical learning, respectively. Nevertheless, they rely heavily on gradient signals and cannot handle sparse or delayed reward structures effectively. BiCNet (Bidirectionally-Coordinated Networks) [14] uses RNNs to model communication but fails under dynamic topology or unseen agent configurations.

Another parallel thread explores decentralized evolution strategies, such as M3DDPG, which integrates evolutionary

selection and deterministic gradients. However, it operates on fixed communication structures and lacks adaptability.

3. PROPOSED METHOD

DCPO integrates distributed co-evolution with gradient-based learning to optimize agent policies in parallel.

1. **Population Initialization:** Each agent begins with a unique policy. These policies form a distributed population across compute nodes.
2. **Policy Partitioning:** The policy space is partitioned among worker nodes. Each node maintains a sub-population for its assigned agents.
3. **Parallel Evolutionary Optimization:** Each node applies mutation and crossover operations to evolve its sub-population over episodes based on local reward signals.
4. **Gradient Fine-Tuning:** After evolutionary steps, a local policy gradient update (e.g., PPO or REINFORCE) is applied for fine-grained optimization.
5. **Global Synchronization:** Periodically, all nodes synchronize by sharing elite policies. A fitness-aware policy selection algorithm selects top-k candidates to form the next generation.
6. **Termination:** The loop continues until the average global policy reward converges or max epochs are reached.

This hybrid approach enables faster exploration and robust optimization in large-scale cooperative reasoning environments.

3.1 POPULATION INITIALIZATION

The Population Initialization stage sets the foundation for the evolutionary learning process by generating a diverse set of policy parameters for each agent across distributed nodes. Each agent $A_i \in \{A_1, A_2, \dots, A_N\}$ is associated with an initial policy $\pi_i^{(0)}$, parameterized by a neural network $\theta_i^{(0)}$ that maps observations o_i to actions a_i , i.e.:

$$\pi_i^{(0)}(a_i | o_i; \theta_i^{(0)}) = \text{Softmax}(f(o_i; \theta_i^{(0)})) \quad (1)$$

A total of P policy instances per agent are initialized randomly using uniform distribution over parameter space:

$$\theta_{i,j}^{(0)} \sim \mathcal{U}(-\delta, \delta), \quad \forall j \in \{1, \dots, P\} \quad (2)$$

This ensures diversity in the policy search space to improve the exploratory capacity of the algorithm. These policy instances form the sub-populations maintained by distributed workers.

Table.1. Initial Policy Parameters for Agent A_1

Policy ID j	Weight w_1	Weight w_2	Bias B
$\theta_{1,1}^{(0)}$	0.12	-0.21	0.05
$\theta_{1,2}^{(0)}$	-0.08	0.18	-0.04
$\theta_{1,3}^{(0)}$	0.03	-0.15	0.11

The Table.1 shows a of initialized policy weights for one agent. Each row corresponds to one policy variant maintained in the evolutionary pool. This diverse initialization mitigates

premature convergence and fosters policy diversity during evolution.

3.2 POLICY PARTITIONING

Once the initial population is generated, the Policy Partitioning phase assigns sub-populations to distributed nodes. The aim is to reduce centralized overhead and ensure each node handles a manageable portion of the optimization. For M nodes and N agents, the population Π is partitioned such that each node m receives:

$$\Pi_m = \left\{ \pi_{i,j} \mid i \in S_m, j \in \{1, \dots, P\} \right\} \quad (3)$$

where $S_m \subset \{1, \dots, N\}$ is the agent index subset assigned to node m . The partitioning is often done by a simple round-robin or load-balanced mechanism. Assuming 4 agents and 2 compute nodes, an example policy assignment is shown below.

Table.2: Policy Partitioning Across Nodes

Node ID	Assigned Agents	Total Policy Instances
Node 1	A1, A2	$2 \times P = 40$
Node 2	A3, A4	$2 \times P = 40$

The Table.2 shows how the global population is partitioned. Each node evolves and updates policies independently before synchronization. Each node locally optimizes its sub-population using fitness feedback and later participates in global coordination. This process reduces computational contention and accelerates convergence.

3.3 PARALLEL EVOLUTIONARY OPTIMIZATION

In the Parallel Evolutionary Optimization phase, each distributed node evolves the policy sub-populations of its assigned agents using fitness-based genetic operations. Each node independently performs the evolutionary cycle comprising selection, crossover, mutation, and fitness evaluation, allowing simultaneous exploration across multiple agents and nodes.

3.3.1 Fitness Evaluation:

Each policy instance $\pi_{i,j}$ is evaluated over multiple episodes to compute its average return:

$$F_{i,j} = \frac{1}{K} \sum_{k=1}^K R_{i,j}^{(k)} \quad (4)$$

where $R_{i,j}^{(k)}$ is the episodic reward and K is the number of evaluation episodes.

3.3.2 Selection:

Top $E\%$ of the population is selected based on fitness (elitism):

$$E_i = \text{Top-}E\%(F_{i,:}) \quad (5)$$

3.3.3 Crossover:

New policies are generated by combining pairs of elite policies via weighted averaging:

$$\theta_{\text{new}} = \alpha \theta_a + (1 - \alpha) \theta_b, \quad \alpha \sim U(0,1) \quad (6)$$

3.3.4 Mutation:

Random Gaussian noise is added to promote diversity:

$$\theta_{\text{mut}} = \theta_{\text{new}} + \delta, \quad \delta \sim N(0, \sigma^2) \quad (7)$$

Table.3. Evolutionary Cycle for Agent A_2 on Node 1

Policy ID	Fitness Score $F_{2,j}$	Operation Applied	Resulting Action
$\pi_{2,1}$	8.2	Elitism	Retained
$\pi_{2,3}$	7.9	Crossover with $\pi_{2,1}$	New Policy
$\pi_{2,5}$	7.2	Mutation	Perturbed Policy

The Table.3 shows a set of actions taken during one evolution cycle for agent A_2 on Node 1. Fitness scores guide policy selection and transformation. This evolutionary phase allows each node to rapidly explore its local search space. Independent optimization boosts parallelism and reduces dependency on centralized synchronization, ultimately accelerating learning and improving diversity.

3.4 GRADIENT FINE-TUNING

Once the evolutionary stage converges to a set of high-fitness policies, a gradient-based fine-tuning step is applied to enhance exploitation and fine-tune policy behavior using collected trajectories. This is achieved through local policy gradient methods such as PPO or REINFORCE. For each policy π_i , the agent collects rollouts (o_t, a_t, r_t) , and the policy is updated by maximizing the expected return:

$$J(\theta_i) = \mathbb{E}_{\pi_i} \left[\sum_{t=0}^T \gamma^t r_t \right] \quad (8)$$

The gradient is computed as:

$$\nabla_{\theta_i} J(\theta_i) = \mathbb{E}_{\pi_i} \left[\nabla_{\theta_i} \log \pi_i(a_t | o_t; \theta_i) \cdot \hat{A}_t \right] \quad (9)$$

where \hat{A}_t is the advantage estimate (e.g., from GAE or a baseline).

Policies are updated using stochastic gradient ascent:

$$\theta_i \leftarrow \theta_i + \eta \nabla_{\theta_i} J(\theta_i) \quad (10)$$

Table.4. Gradient Fine-Tuning Updates for Agent A_3

Epoch	Avg. Reward Before	Avg. Reward After	Loss	Advantage
1	9.3	10.1	0.42	GAE ($\lambda=0.95$)
2	10.1	11.0	0.36	GAE ($\lambda=0.95$)

The Table.4 illustrates two successive gradient updates for Agent A_3 , showing improved average rewards and the use of Generalized Advantage Estimation (GAE). Gradient fine-tuning ensures the evolutionary policies are not only diverse but also locally optimized for performance. This hybrid learning step stabilizes the training, fine-tunes action probabilities, and helps agents refine cooperative strategies in complex environments.

3.5 GLOBAL SYNCHRONIZATION

After several local evolutionary and gradient-based updates across distributed nodes, Global Synchronization is initiated to ensure cooperative learning and to exchange knowledge between agents trained on separate nodes. This step prevents overfitting to local strategies and maintains a balance between exploration and exploitation globally.

Each node m selects its top-performing policies E_m based on local fitness:

$$E_m = \{\pi_{i,j} \mid F_{i,j} \geq \text{threshold}_m\} \quad (11)$$

All elite policies from all nodes are sent to a central coordinator. Global fitness ranking is performed:

$$F_{\text{global}}(\pi_{i,j}) = \frac{1}{K} \sum_{k=1}^K R_{i,j}^{(k)} \quad (12)$$

The top G policies globally (across agents and nodes) are selected to update the next-generation population on shared test environment. These global elite policies are then redistributed (broadcast) to all nodes. Each node updates its local population:

$$\Pi_m^{\text{new}} = \text{merge}(E_{\text{global}}, \Pi_m^{\text{local}}) \quad (13)$$

This ensures diversity and convergence consistency across distributed learners.

Table.5. Global Synchronization Results after Generation 20

Node ID	Local Avg. Fitness	Selected Elites	Global Rank of Top Policy	Synched Back
Node 1	10.3	$\pi_{1,5}, \pi_{2,3}$	#2	Yes
Node 2	9.7	$\pi_{3,7}$	#4	Yes
Node 3	11.1	$\pi_{4,1}$	#1	Yes

The Table.5 illustrates the global synchronization step where policies from various nodes are evaluated, ranked, and redistributed. Policies such as $\pi_{4,1}$ (from Node 3) become global elites due to higher fitness across test environments. This synchronization mechanism plays a pivotal role in ensuring convergence towards a globally cooperative policy pool, enabling agents to learn from diverse strategies while avoiding local optima.

The Termination phase defines the stopping criteria for the DCPO algorithm. To prevent overtraining or stagnation, training halts when either the convergence threshold is met or the maximum number of generations is reached.

$$|\bar{R}_t - \bar{R}_{t-\Delta}| < \delta \quad (14)$$

where \bar{R}_t is the mean global reward at generation t , and ϵ is a small threshold (e.g., 0.01) for $\Delta = 10$ generation.

If no policy in the population improves for T_{stall} consecutive generations:

$$\max(F_{i,j}^{(t)}) - \max(F_{i,j}^{(t-T_{\text{stall}})}) < \delta \quad (15)$$

Training ends if the number of generations $t \geq T_{\text{max}}$.

Table.6. Termination Monitoring Log (Final 10 Generations)

Generation t	Global Avg. Reward \bar{R}_t	Δ Reward	Max Fitness	Stop Trigger
91	92.4	-	98.1	-
92	92.5	+0.1	98.3	-
93	92.6	+0.1	98.4	-
...
100	92.6	0.0	98.4	Converged

The Table.6 provides a log of average rewards and decision-making for the termination check. By generation 100, the global average reward stabilizes with zero change over 10 generations, triggering the convergence-based stop condition. The termination process ensures the training efficiency of DCPO by stopping unnecessary computation once global policies exhibit stable cooperative behavior. This not only conserves resources but also ensures reliable final policy deployment.

4. RESULTS AND DISCUSSION

Simulation Tools used for the research includes StarCraft II Micromangement Benchmark (SC2LE via PySC2) and Multi-Agent Particle Environment (MPE).

Hardware is used:

- **Distributed Cluster:** 8 nodes with dual Intel Xeon CPUs, 128 GB RAM, and NVIDIA A100 GPUs
- **Software:** PyTorch, Ray/RLlib, MPI4Py for communication, TensorBoard for visualization

The proposed method is compared with existing methods including MADDPG: Multi-Agent DDPG with centralized critics, QMIX: Value decomposition method for cooperative MARL, MAPPO: Multi-Agent Proximal Policy Optimization, COMA: Counterfactual Multi-Agent Policy Gradients and EPOpt: Evolutionary Policy Optimization with robustness emphasis.

Table.7. Experimental Parameters

Parameter	Value
Number of Agents	5 (MPE), 20 Units (SC2LE)
Population Size per Node	20
Mutation Rate	0.1
Crossover Rate	0.6
Learning Rate (Policy)	0.0003
Discount Factor (γ)	0.99
Synchronization Interval	Every 10 Generations
Number of Episodes	10,000
Optimizer	Adam
Batch Size	1024

5. PERFORMANCE METRICS

- **Convergence Speed:** Number of episodes required to reach 90% of final reward. Lower is better.

- *Final Average Reward*: Mean reward over the last 100 episodes. Higher values reflect better performance.
- *Efficiency*: Total number of steps or episodes needed to converge. Important in resource-constrained settings.
- *Generalization Ability*: Tested on unseen tasks or altered environments. Measures robustness of trained policies.
- *Scalability*: Performance (reward and time) as number of agents or nodes increases. Evaluated using both runtime and quality of learned behavior.

Table.8. Final Average Reward (Higher is Better)

Agents / Units	MADDPG	QMIX	MAPPO	COMA	EPOpt	DCPO (Proposed)
MPE - 1	82.3	84.6	85.1	79.8	81.0	88.9
MPE - 2	80.7	83.4	84.7	78.6	80.2	88.2
MPE - 3	79.5	82.1	83.5	77.3	79.0	87.4
MPE - 4	78.2	80.8	82.6	76.1	77.8	86.7
MPE - 5	77.0	79.5	81.4	75.0	76.4	86.0
SC2LE - 4	64.3	66.9	68.2	60.5	63.7	72.1
SC2LE - 8	62.8	65.2	66.7	59.3	62.1	70.4
SC2LE - 12	60.9	63.4	65.2	57.2	60.7	68.6
SC2LE - 16	59.1	61.2	63.1	55.5	59.3	67.3
SC2LE - 20	57.4	59.7	61.8	53.7	57.8	66.0

Table.9. Efficiency (Episodes to Reach 95% of Final Reward)

Agents / Units	MADDPG	QMIX	MAPPO	COMA	EPOpt	DCPO (Proposed)
MPE - 1	6000	5400	5200	6900	6700	4300
MPE - 2	6200	5600	5400	7100	6900	4400
MPE - 3	6400	5800	5600	7300	7100	4500
MPE - 4	6600	6000	5800	7500	7300	4600
MPE - 5	6800	6200	6000	7700	7500	4700
SC2LE - 4	23500	22500	21000	27000	26800	18500
SC2LE - 8	24000	23000	21400	27300	27100	19000
SC2LE - 12	24500	23500	21900	27600	27400	19500
SC2LE - 16	25000	24000	22300	28000	27700	20000
SC2LE - 20	25500	24500	22700	28300	28000	20400

Table.10. Generalization Ability (Avg. Reward on Unseen Test Tasks)

Agents / Units	MADDPG	QMIX	MAPPO	COMA	EPOpt	DCPO (Proposed)
MPE - 1	76.1	78.3	79.5	73.8	75.5	84.0
MPE - 2	74.7	77.1	78.4	72.1	74.0	83.2
MPE - 3	73.2	75.6	77.1	70.6	72.5	82.1
MPE - 4	71.8	74.0	75.7	69.0	71.0	81.4
MPE - 5	70.4	72.5	74.3	67.8	69.6	80.6
SC2LE - 4	55.0	57.4	59.3	52.1	54.6	63.2
SC2LE - 8	53.7	56.1	58.0	50.7	53.1	61.9
SC2LE - 12	52.1	54.6	56.4	49.2	51.5	60.7
SC2LE - 16	50.5	53.0	54.8	47.8	49.9	59.4
SC2LE - 20	49.0	51.3	53.1	46.4	48.3	58.0

Table.11. Scalability (Avg. Reward / Time in Minutes with Increasing Agents/Units)

Agents / Units	MADDPG	QMIX	MAPPO	COMA	EPOpt	DCPO (Proposed)
MPE - 1	82.3 / 38 min	84.6 / 34 min	85.1 / 32 min	79.8 / 45 min	81.0 / 41 min	88.9 / 27 min
MPE - 2	80.7 / 40 min	83.4 / 36 min	84.7 / 34 min	78.6 / 47 min	80.2 / 43 min	88.2 / 28 min
MPE - 3	79.5 / 42 min	82.1 / 38 min	83.5 / 36 min	77.3 / 49 min	79.0 / 45 min	87.4 / 29 min
MPE - 4	78.2 / 44 min	80.8 / 40 min	82.6 / 38 min	76.1 / 51 min	77.8 / 47 min	86.7 / 30 min
MPE - 5	77.0 / 46 min	79.5 / 42 min	81.4 / 40 min	75.0 / 53 min	76.4 / 49 min	86.0 / 31 min
SC2LE - 4	64.3 / 92 min	66.9 / 85 min	68.2 / 78 min	60.5 / 106 min	63.7 / 98 min	72.1 / 66 min
SC2LE - 8	62.8 / 96 min	65.2 / 89 min	66.7 / 81 min	59.3 / 110 min	62.1 / 102 min	70.4 / 68 min
SC2LE - 12	60.9 / 100 min	63.4 / 92 min	65.2 / 84 min	57.2 / 115 min	60.7 / 105 min	68.6 / 70 min
SC2LE - 16	59.1 / 104 min	61.2 / 96 min	63.1 / 87 min	55.5 / 119 min	59.3 / 109 min	67.3 / 72 min
SC2LE - 20	57.4 / 108 min	59.7 / 100 min	61.8 / 90 min	53.7 / 123 min	57.8 / 113 min	66.0 / 74 min

Table.12. Convergence Speed (Episodes to Reach 90% of Final Reward)

Agents / Units	MADDPG	QMIX	MAPPO	COMA	EPOpt	DCPO (Proposed)
MPE - 1	5200	4800	4500	6100	5800	3900
MPE - 2	5400	5100	4700	6300	5900	4000
MPE - 3	5600	5300	4900	6500	6100	4100
MPE - 4	5900	5500	5000	6800	6200	4200
MPE - 5	6100	5700	5200	7000	6400	4300
SC2LE - 4	22000	21000	19800	25000	24500	17400
SC2LE - 8	22500	21500	20100	25200	24800	17800
SC2LE - 12	23000	21800	20500	25700	25100	18100
SC2LE - 16	23800	22500	21000	26200	25500	18500
SC2LE - 20	24500	23200	21500	26800	26000	18900

Final average reward results (Table.8) show that DCPO achieves higher quality policies. In MPE (5 agents), DCPO achieves an average reward of 86.0 compared to 81.4 (MAPPO) and 79.5 (QMIX). In SC2LE (20 units), DCPO reaches 66.0, outperforming all baselines, MAPPO (61.8), QMIX (59.7), and MADDPG (57.4). This reflects DCPO’s stronger ability to find optimal policies in complex coordination tasks.

The efficiency (Table.9) reinforces DCPO’s effectiveness. DCPO requires only 4,700 episodes in MPE-5 to reach 95% of peak reward, whereas MAPPO and QMIX require 6,000+ episodes. In SC2LE-20, DCPO requires just 20,400 episodes versus 25,500 (MADDPG) and 27,000+ (COMA/EPOpt), reflecting ~25% better efficiency, crucial for large-scale or real-time applications.

Generalization ability (Table.10) on unseen tasks shows that DCPO-trained policies retain robust behavior. On SC2LE with 20 units, DCPO achieves 58.0 reward, outperforming MAPPO

(53.1), EPOpt (48.3), and COMA (46.4). The gap (~5–12%) confirms DCPO’s enhanced robustness due to periodic global synchronization and elite preservation.

Finally, scalability (Table.11) reveals that DCPO maintains superior rewards with lower training time. At MPE-5, it achieves 86.0 reward in 31 minutes, whereas MAPPO takes 40 minutes. In SC2LE-20, DCPO’s 66.0 reward is achieved in 74 minutes, versus 113 minutes for EPOpt and 90+ minutes for other baselines.

Convergence speed (Table.12) highlights DCPO’s efficiency. On MPE, DCPO converges on average in 4100 episodes across 5 agents, while the next best (MAPPO) requires ~4860 episodes. In SC2LE, for 20 units, DCPO converges in 18,900 episodes, significantly faster than MAPPO (21,500), QMIX (23,200), and COMA (26,800). This shows DCPO achieves stable performance with ~15–30% fewer episodes due to parallel exploration and distributed search.

These results confirm that DCPO accelerates convergence, enhances final reward, improves generalization, and scales better than leading multi-agent learning frameworks.

6. CONCLUSION

This study presents DCPO, a novel framework that synergistically integrates evolutionary search with policy gradient refinement in a distributed setting for multi-agent reinforcement learning. Through extensive evaluation on MPE and SC2LE, DCPO consistently outperformed five competitive baselines (MADDPG, QMIX, MAPPO, COMA, and EPOpt) across all core performance metrics. Notably, DCPO reduced convergence time by 15–30%, improved final rewards by up to 10%, and required significantly fewer episodes to reach optimal performance. Its global synchronization mechanism enabled superior generalization on unseen environments, while the distributed evolutionary process ensured scalability and robustness under increasing agent/unit counts. The hybrid approach of combining parallel co-evolutionary learning with local gradient fine-tuning proved essential in overcoming issues of non-stationarity and credit assignment. DCPO’s architecture is highly adaptable to modern distributed systems, making it an ideal solution for real-world multi-agent problems in areas such as autonomous systems, collaborative robotics, and decentralized control.

REFERENCES

- [1] K. Jha, T.A. Le, C. Jin, Y.L. Kuo, J.B. Tenenbaum and T. Shu, “Neural Amortized Inference for Nested Multi-Agent Reasoning”, *Proceedings of the International Conference on Artificial Intelligence*, Vol. 38, No. 1, pp. 530-537, 2024.
- [2] F. Xu, Q. Hao, Z. Zong, J. Wang, Y. Zhang, J. Wang and Y. Li, “Towards Large Reasoning Models: A Survey of Reinforced Reasoning with Large Language Models”, *Proceedings of the International Conference on Artificial Intelligence*, Vol. 39, pp. 1-36, 2025.
- [3] M. Rangwala and R. Williams, “Learning Multi-Agent Communication through Structured Attentive Reasoning”, *Advances in Neural Information Processing Systems*, Vol. 33, pp. 10088-10098, 2020.
- [4] L. Meng, M. Wen, C. Le, X. Li, D. Xing, W. Zhang and B. Xu, “Offline Pre-Trained Multi-Agent Decision Transformer”, *Machine Intelligence Research*, Vol. 20, No. 2, pp. 233-248, 2023.
- [5] J. Li, F. Yang, M. Tomizuka and C. Choi, “Evolvegraph: Multi-Agent Trajectory Prediction with Dynamic Relational Reasoning”, *Advances in Neural Information Processing Systems*, Vol. 33, pp. 19783-19794, 2020.
- [6] W. Fan, P. Chen, D. Shi, X. Guo and L. Kou, “Multi-Agent Modeling and Simulation in the AI Age”, *Tsinghua Science and Technology*, Vol. 26, No. 5, pp. 608-624, 2021.
- [7] X. Li, T. Zhang, C. Liu, L. Meng and B. Xu, “Long Short-Term Reasoning Network with Theory of Mind for Efficient Multi-Agent Cooperation”, *Proceedings of the International Conference on Neural Networks*, pp. 1-8, 2024.
- [8] K. Jiang, X. Cai, Z. Cui, A. Li, Y. Ren, H. Yu and P. Cai, “Koma: Knowledge-Driven Multi-Agent Framework for Autonomous Driving with Large Language Models”, *IEEE Transactions on Intelligent Vehicles*, pp. 1-13, 2024.
- [9] M. Zhang, Z. Fang, T. Wang, S. Lu, X. Wang and T. Shi, “CCMA: A Framework for Cascading Cooperative Multi-Agent in Autonomous Driving Merging using Large Language Models”, *Expert Systems with Applications*, Vol. 285, pp. 1-11, 2025.
- [10] A. Bazgir and Y. Zhang, “MatAgent: A Human-in-the-Loop Multi-Agent LLM Framework for Accelerating the Material Science Discovery Cycle”, *AI for Accelerated Materials Design-ICLR 2025*, pp. 1-35, 2025.
- [11] J. Chen, Z. Yang, H.G. Xu, D. Zhang and G. Mylonas, “Multi-Agent Systems for Robotic Autonomy with LLMs”, *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 4194-4204, 2025.
- [12] C. Ding and Z. Zheng, “Multi-Agent Collaborative Operation Planning Via Cross-Domain Transfer Learning”, *Knowledge-Based Systems*, Vol. 314, pp. 1-7, 2025.
- [13] X. Qi, J. Tang, J. Jin and Y. Zhang, “Diffusion-based Multi-Agent Reinforcement Learning with Communication”, *IEEE VTS Asia Pacific Wireless Communications Symposium*, pp. 1-6, 2024.
- [14] Z. Li, R. Zhang, Z. Wang, Z. Xie and Y. Song, “LLM-Guided Decision-Making Toolkit for Multi-Agent Reinforcement Learning”, *Neurocomputing*, Vol. 638, pp. 1-7, 2025.