

# MisoDICE: Multi-Agent Imitation from Unlabeled Mixed-Quality Demonstrations

- The Viet Bui, Singapore Management University, Singapore
- Tien Anh Mai, Singapore Management University, Singapore
- Thanh Hong Nguyen, University of Oregon, USA



## Introduction

- **The Problem:** Obtaining high-quality expert data in multi-agent environments is expensive and impractical due to complex joint state-action spaces. Real-world data is often **unlabeled and mixed-quality**, containing both expert and sub-optimal trajectories.
- **The Challenge:** Existing methods assume access to expert labels or high-quality demonstrations. Learning from unlabeled mixed data requires distinguishing expert behaviors without ground-truth rewards.

## Preliminaries

**Setting:** Cooperative MARL modeled as a POMDP:

$$M = \langle S, A, P, r, Z, O, n, N, \gamma \rangle$$

**Data:** We operate on an Unlabeled Dataset ( $\mathcal{D}_{\text{unlabeled}}$ ) containing a mix of expert and non-expert trajectories.

**Goal:** Recover optimal local policies  $\pi_i$  without access to the ground-truth reward function, maximizing the expected return based on the inferred expert distribution.

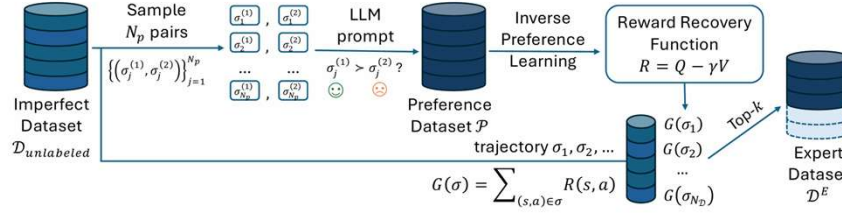
## Experiments & Results

|            | BC              |                 |                 | OMAPL      | INDD       | MARL-SL    | VDN        | MisoDICE (ours)   |
|------------|-----------------|-----------------|-----------------|------------|------------|------------|------------|-------------------|
|            | $(\beta = 0.0)$ | $(\beta = 0.5)$ | $(\beta = 1.0)$ |            |            |            |            |                   |
| 2c_vs_64zg | 8.5 ± 0.1       | 9.7 ± 0.3       | 12.6 ± 0.3      | 12.2 ± 0.4 | 14.6 ± 1.0 | 14.0 ± 1.6 | 12.7 ± 0.6 | <b>16.4 ± 1.3</b> |
| 5m_vs_6m   | 5.0 ± 1.1       | 6.7 ± 0.0       | 6.1 ± 0.1       | 5.7 ± 0.2  | 6.7 ± 0.1  | 6.8 ± 0.1  | 6.2 ± 1.4  | <b>7.3 ± 0.1</b>  |
| 6h_vs_8z   | 7.0 ± 0.0       | 7.4 ± 0.0       | 7.2 ± 0.1       | 6.6 ± 0.2  | 7.5 ± 0.2  | 7.8 ± 0.1  | 8.2 ± 0.2  | <b>8.7 ± 0.2</b>  |
| corridor   | 1.5 ± 0.1       | 1.5 ± 0.2       | 4.3 ± 0.7       | 2.2 ± 1.3  | 4.4 ± 1.2  | 1.8 ± 0.2  | 4.7 ± 0.6  | <b>5.8 ± 0.8</b>  |
| 5_vs_5     | 9.2 ± 0.1       | 11.7 ± 0.5      | 10.2 ± 0.5      | 9.6 ± 1.1  | 10.9 ± 0.1 | 11.6 ± 0.3 | 11.5 ± 0.2 | <b>12.4 ± 0.5</b> |
| 10_vs_10   | 10.3 ± 0.6      | 11.8 ± 0.5      | 10.6 ± 0.2      | 10.1 ± 0.9 | 11.0 ± 0.7 | 11.9 ± 0.4 | 12.4 ± 0.2 | <b>12.9 ± 0.2</b> |
| 10_vs_11   | 8.2 ± 0.4       | 9.6 ± 0.4       | 8.7 ± 0.3       | 8.5 ± 1.2  | 9.4 ± 0.4  | 9.9 ± 0.3  | 10.4 ± 0.1 | <b>10.7 ± 0.4</b> |
| 20_vs_20   | 10.1 ± 0.2      | 10.4 ± 0.5      | 10.5 ± 0.3      | 9.4 ± 0.4  | 11.4 ± 0.5 | 13.1 ± 0.4 | 12.1 ± 0.5 | <b>13.5 ± 0.5</b> |
| 20_vs_23   | 8.1 ± 0.2       | 8.6 ± 0.3       | 8.3 ± 0.2       | 7.9 ± 0.3  | 9.6 ± 0.3  | 9.6 ± 0.3  | 10.3 ± 0.4 | <b>10.6 ± 0.2</b> |

**Conclusion:** MisoDICE significantly outperforms baselines by effectively leveraging unlabeled mixed-quality data, confirming the benefits of the two-stage labeling and convex value-decomposition approach.

## Phase 1: Expert Identification (Data Labeling)

- Step 1 (LLM Preferences):** Sample trajectory pairs and use an LLM (e.g., GPT-4o) to generate preference labels based on semantic game features (health, position).
- Step 2 (Reward Recovery):** Train O-MAPL (Preference-based MARL) on these labels to learn a soft Q-function. Recover rewards via  $R \approx Q - \gamma V$ .
- Step 3 (Ranking):** Rank trajectories by total recovered return. The top- $k$  are selected as the Expert Dataset ( $\mathcal{D}^E$ ); the rest form the suboptimal set ( $\mathcal{D}^{Mix}$ ).



## Value Factorization & consistency

**Value Factorization:** To handle the combinatorial action space, we decompose the global value function using a linear mixing network to preserve convexity:

$$v^{tot}(s) = \mathcal{M}_\phi[\{v_i(s_i)\}] = \sum \phi_i v_i(s_i) + \phi_0$$

*Note: Non-linear mixing (like ReLU networks) destroys the convexity of the*

*DICE objective, leading to instability.*

**Occupancy Ratio Estimation:** We estimate the density ratio  $w(s, a)$  using a discriminator trained with a linear mixing network to ensure concavity.

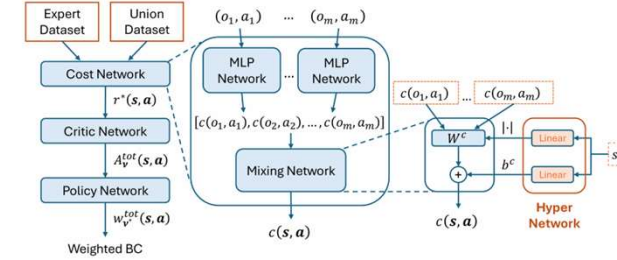
## Phase 2: MisoDICE Algorithm

**Objective:** Minimize the divergence between the learned policy and the expert distribution, regularized by the union (mixed) distribution:

$$\max_{\pi_{tot}} D_{KL}(\rho_{tot}^\pi || \rho_{tot}^E) + \alpha D_{KL}(\rho_{tot}^\pi || \rho_{tot}^U)$$

- $\rho_{tot}^E$ : Expert distribution (from Phase 1).
- $\rho_{tot}^U$ : Union distribution (Expert + Mixed).
- $\alpha$ : Hyperparameter controlling influence of suboptimal data.

**Optimization:** We reformulate this as a convex optimization problem over stationary distributions using the DICE framework.



## Policy Extraction (Global-Local Consistency)

The optimal local policies are recovered via *Weighted Behavior Cloning (WBC)*, ensuring that maximizing local objectives results in the global optimum:

$$\pi_i^*(a_i | s_i) \propto \exp \left( \frac{\phi_i^*}{1 + \alpha} q_i^*(s_i, a_i) + \log \mu_i^U(a_i | s_i) \right)$$

