

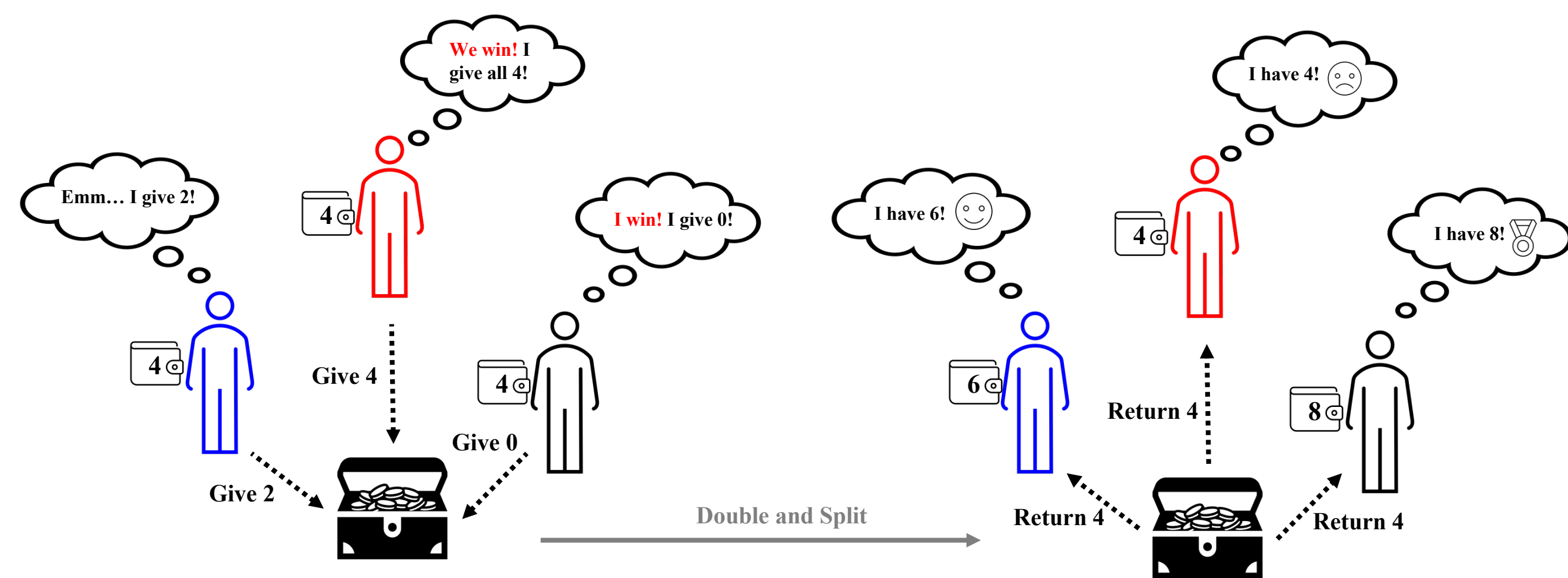
TL;DR

We propose the Altruistic Gradient Adjustment (AgA) for mixed-motive games—a method that aligns individual and collective objectives through gradient adjustments, supported by theoretical proofs and empirical validation.

Background & Overview

A **mixed-motive game** is a game where players have both cooperative and competitive motivations, and their incentives are partly aligned and partly opposed.

Here is an example of mixed-motive game, each player has 4 dollars and need to contribute any dollars to the public pool. The public pool will double the money in pool and split them among players.



In the example, the altruistic player ends up the poorest; the selfish player becomes the wealthiest; the collective outcome is suboptimal.

In this paper, we aim to align individual and collective objectives to foster cooperation in mixed-motive settings.

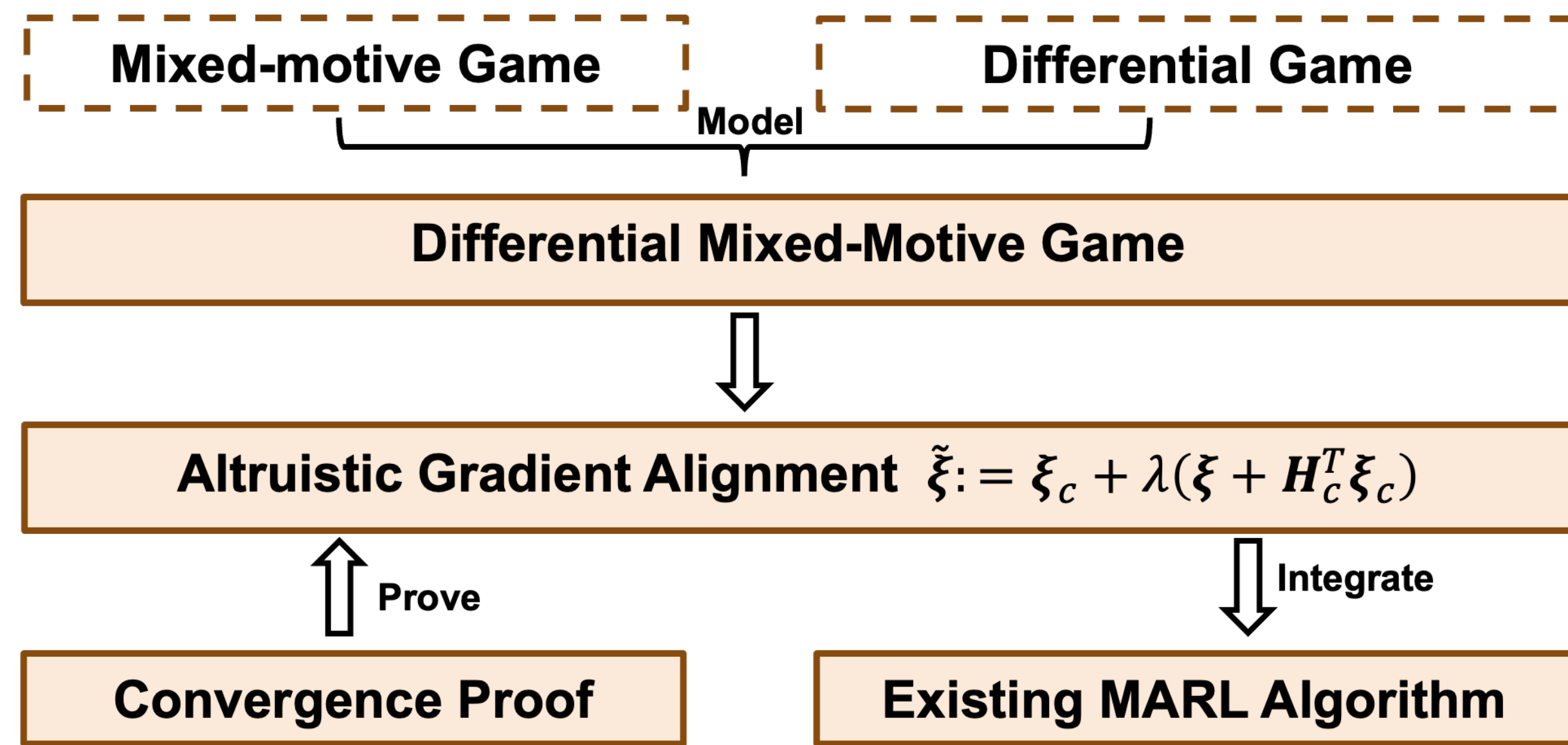


Figure 1. Overview of our proposed altruistic gradient alignment method.

Differential Mixed-motive Game

Differential Mixed-motive Game (DMG) is defined as a tuple $\{\mathcal{N}, \mathbf{w}, \ell\}$, where

- $\mathcal{N} = \{1, \dots, n\}$ is the set of players.
- The parameter set $\mathbf{w} = [\mathbf{w}_i]^n \in \mathbb{R}^d$ is defined, each with $\mathbf{w}_i \in \mathbb{R}^{d_i}$ and $d = \sum_{i=1}^n d_i$.
- $\ell = \{\ell_i : \mathbb{R}^d \rightarrow \mathbb{R}\}_{i=1}^n$ represents the corresponding losses. These losses are assumed to be at least twice differential.

Here, Each player $i \in \mathcal{N}$ is equipped with a policy, parameterized by \mathbf{w}_i , aiming to minimize its loss ℓ_i .

Differentiable losses exhibit the **mixed motivation property**: minimization of individual losses can result in a conflict between individuals or between individual and collective objectives (e.g., maximizing individual stats and winning the game are often conflict in basketball matches).

Altruistic Gradient Adjustment (AgA)

Building on differential mixed-motive game and gradient-adjustment optimization, we propose the **Altruistic Gradient Adjustment (AgA) algorithm to align individual and collective objectives from a gradient-based perspective.**

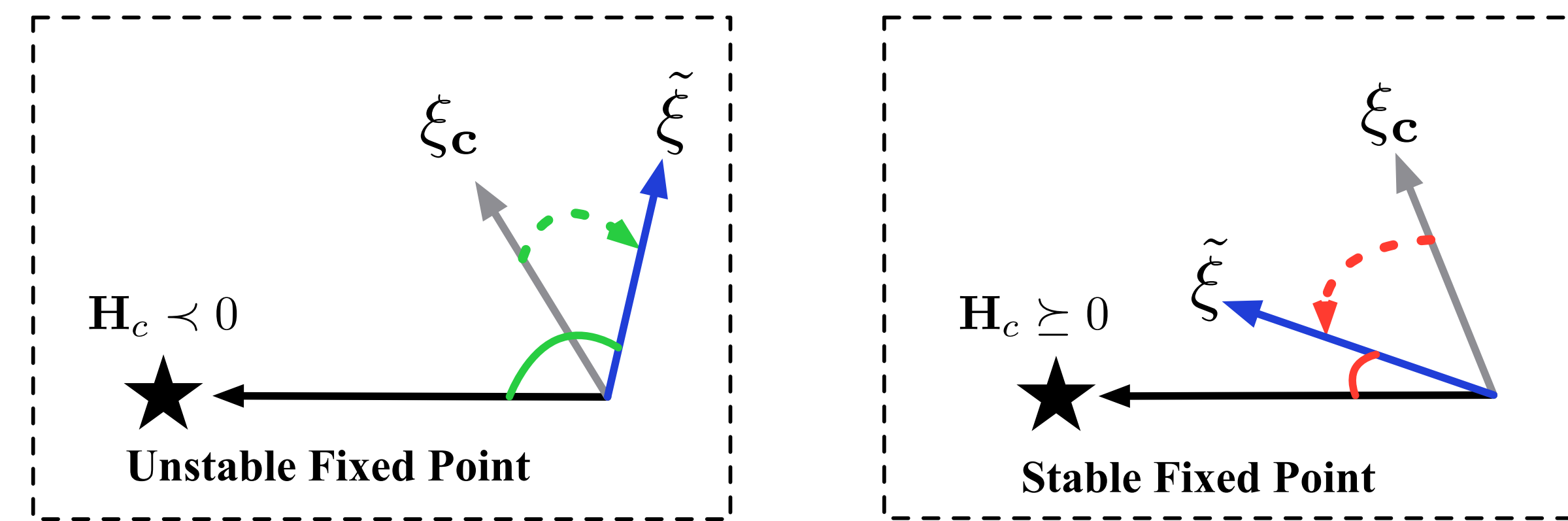
Definition (Altruistic Gradient Adjustment.) Altruistic gradient adjustment (AgA) extends the gradient term in the learning dynamic as

$$\tilde{\xi} := \xi_c + \lambda \xi_{adj} = \xi_c + \lambda (\xi + \mathbf{H}_c^T \xi_c), \quad (1)$$

where $\lambda \in \mathbb{R}$ is alignment parameter, ξ is the gradient of the losses with respect to the parameters of the respective players, $\lambda \xi_{adj}$ is called adjustment term. In ξ_{adj} , ξ_c and \mathbf{H}_c is the gradient vector and Hessian matrix of the game about collective loss.

In our AgA method, we only need to compute the Hessian-vector product $\mathbf{H}_c^T \xi_c$, instead of the full Hessian matrix \mathbf{H}_c^T , reducing the time complexity to $\mathbf{H}_c^T \xi_c$ is $\mathcal{O}(n)$ for n weights.

We theoretically show that selecting the appropriate sign of λ ensures that: (1) AgA pushes the gradient away from unstable points, and (2) AgA pulls the gradient toward stable fixed points via an adjustment term.



Case 1: Push out of Unstable Fixed Point

Case 2: Pull toward Stable Fixed Point

Corollary. In the neighborhood of fixed points of the collective objective, AgA will pull the gradient toward stable fixed points, which means $\theta(\tilde{\xi}, \nabla \mathcal{H}_c) \leq \theta(\xi_c, \nabla \mathcal{H}_c)$, and push away from unstable ones, indicated by $\theta(\tilde{\xi}, \nabla \mathcal{H}_c) \geq \theta(\xi_c, \nabla \mathcal{H}_c)$, if λ satisfies $\lambda \cdot \langle \xi_c, \nabla \mathcal{H}_c \rangle (\langle \xi, \nabla \mathcal{H}_c \rangle + \|\nabla \mathcal{H}_c\|^2) \geq 0$.

Case Study: A Toy Experiment

We conducted a series of toy experiments in a two-player differentiable mixed-motive game. The game is defined as follows.

Example. Consider a two-player DMG with $\ell_1(a_1, a_2) = -\sin(a_1 a_2 + a_2^2)$ and $\ell_2(a_1, a_2) = -[\cos(1 + a_1 - (1 + a_2)^2) + a_1 a_2^2]$, where a_i represents the action of the player i ($i = 1, 2$), and $a_i \in \mathbb{R}$. The rewards for the two players are the negation of their respective losses.

AgA Successfully Aligning Individual and Collective Objectives

Fig. 2a shows the trajectories over the collective reward landscape, with deeper orange indicating higher rewards. Only Simul-Co and AgA move toward the social optimum. However, Simul-Co ignores Player 1's interests, focusing on its own reward peaks and valleys.

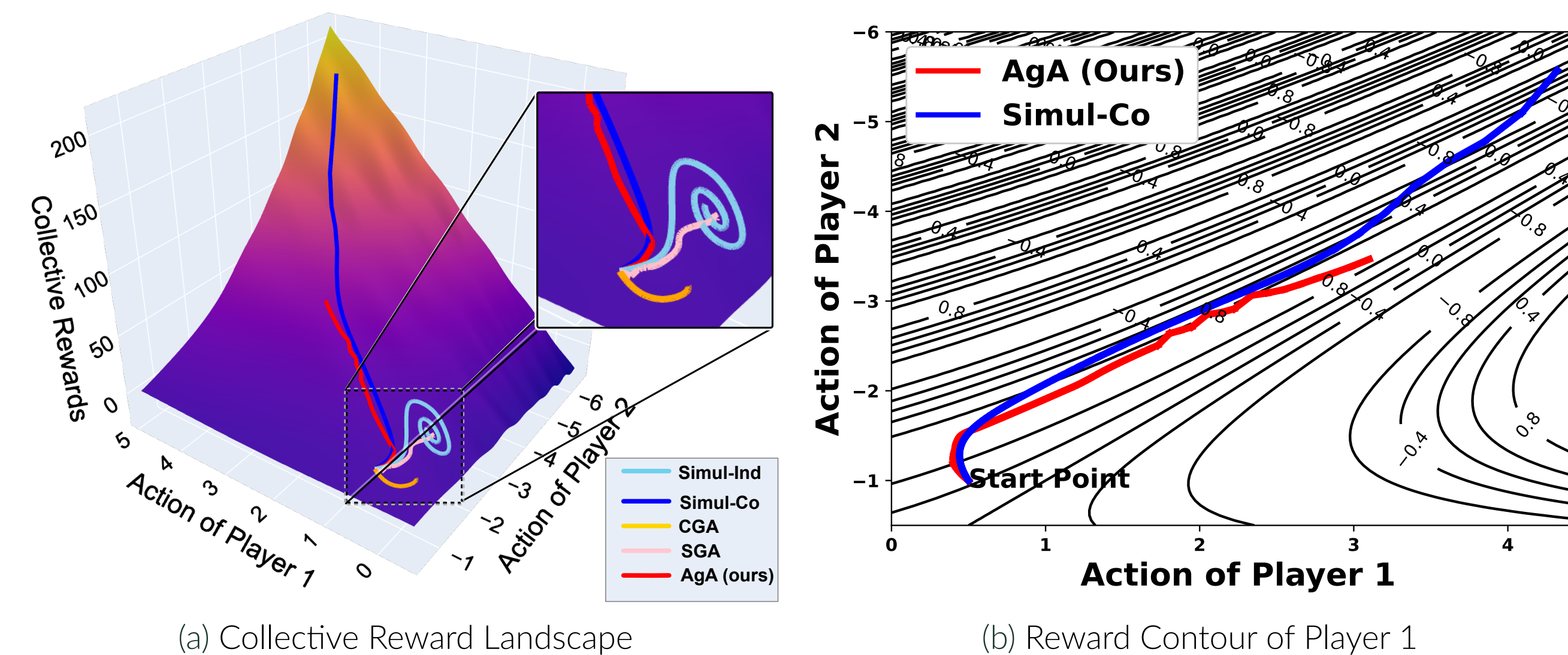
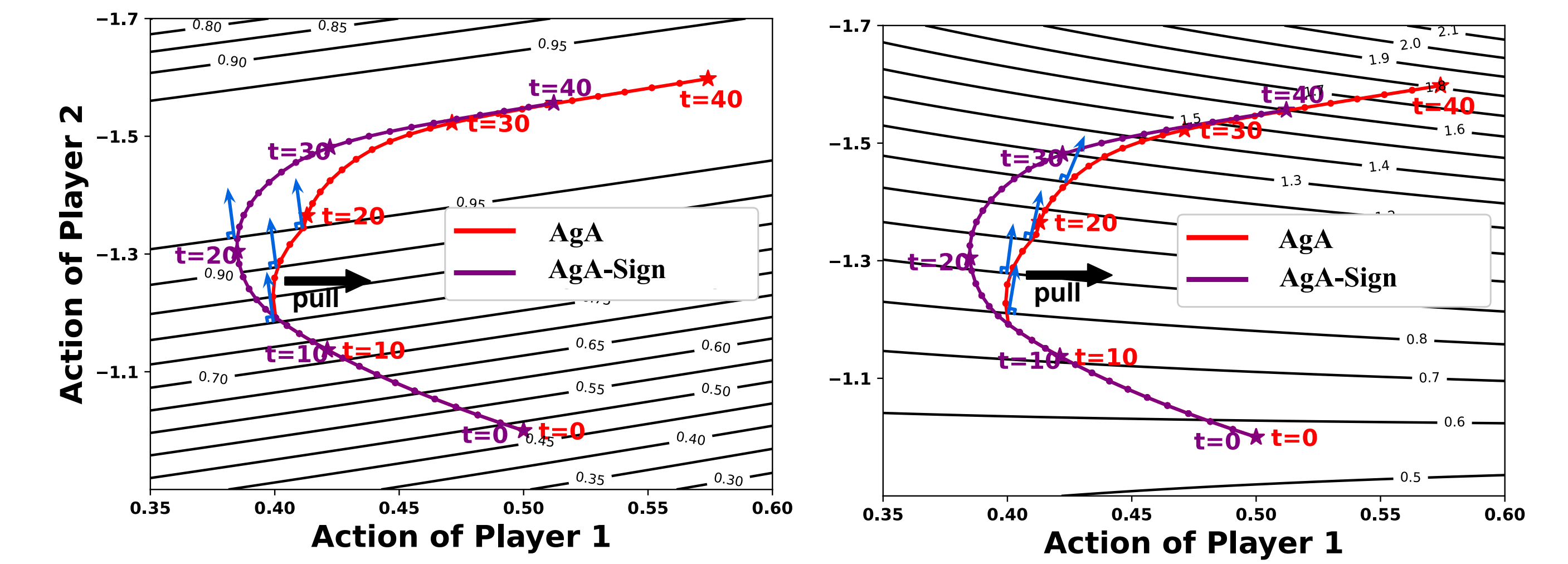


Figure 2. Trajectories of optimization in a two-player DMG.

Alignment Effectiveness of Corollary 4.3

The comparison between AgA (shown in red) and AgA without sign alignment (AgA-Sign, in purple) trajectories spans 40 steps, marked at every tenth step. Starting from the 14th step, sign alignment pulls the gradient toward the steepest direction, resulting in AgA reducing the number of steps by approximately 15% compared to AgA-Sign by the end of the trajectory.



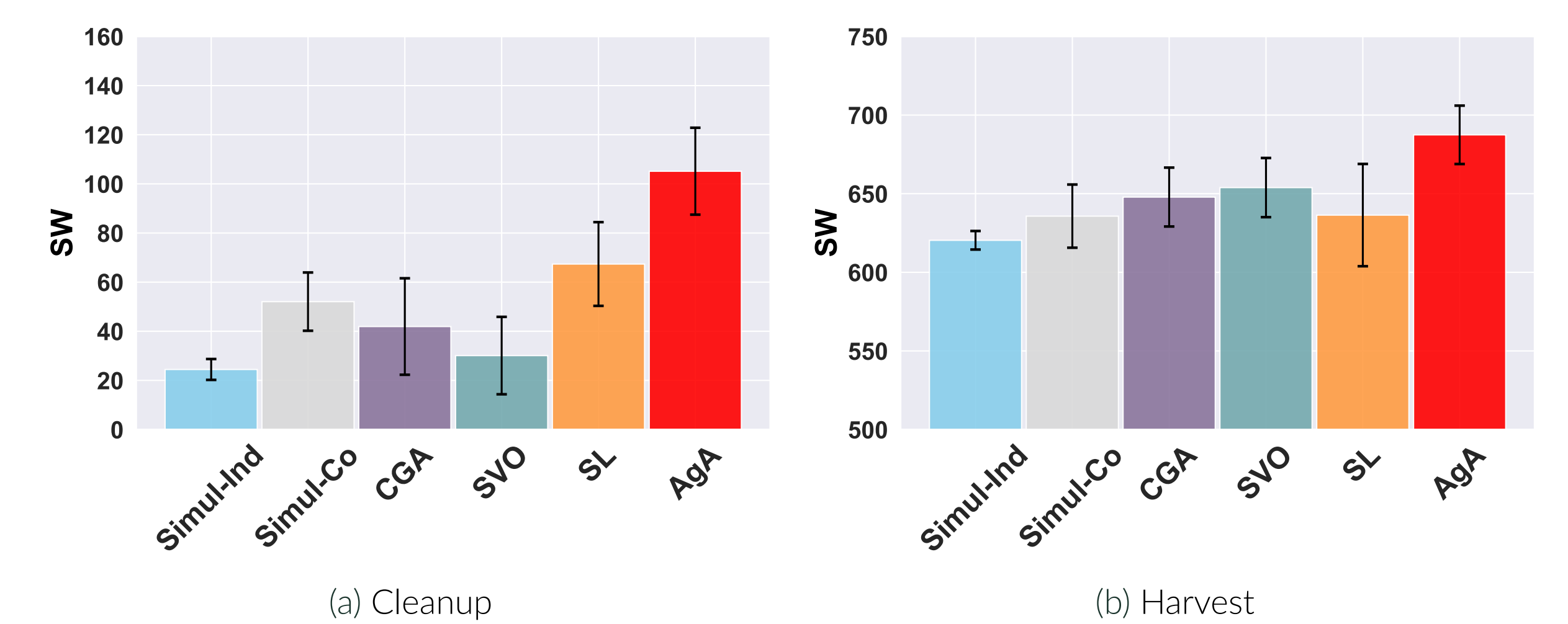
Main Experiment Results

Two-player public good games

A two-player public goods matrix game: Players 1, 2 contribute amounts a_i from a budget $[0, b]$, and their payoffs are $p_i = b - a_i + \frac{c}{2}(a_1 + a_2)$. Social welfare is $SW = p_1 + p_2$. In our experiments, we set $b = 1$ and $c = 1.5$. We show the mean of value and 95% confidence interval utilizing 50 random runs.

Metrics	Simul-Ind	CGA	SGA	SVO	Simul-Co	SL	AgA
r_1	1.133 ± 0.063	1.156 ± 0.060	1.175 ± 0.062	1.104 ± 0.054	1.433 ± 0.056	1.314 ± 0.062	1.443 ± 0.042
r_2	1.184 ± 0.065	1.150 ± 0.057	1.137 ± 0.063	1.060 ± 0.051	1.381 ± 0.065	1.371 ± 0.057	1.459 ± 0.041
SW	2.316 ± 0.039	2.306 ± 0.039	2.312 ± 0.044	2.164 ± 0.026	2.814 ± 0.033	2.684 ± 0.049	2.903 ± 0.023
Equality	0.923 ± 0.014	0.929 ± 0.012	0.924 ± 0.013	0.930 ± 0.011	0.941 ± 0.014	0.940 ± 0.011	0.960 ± 0.008

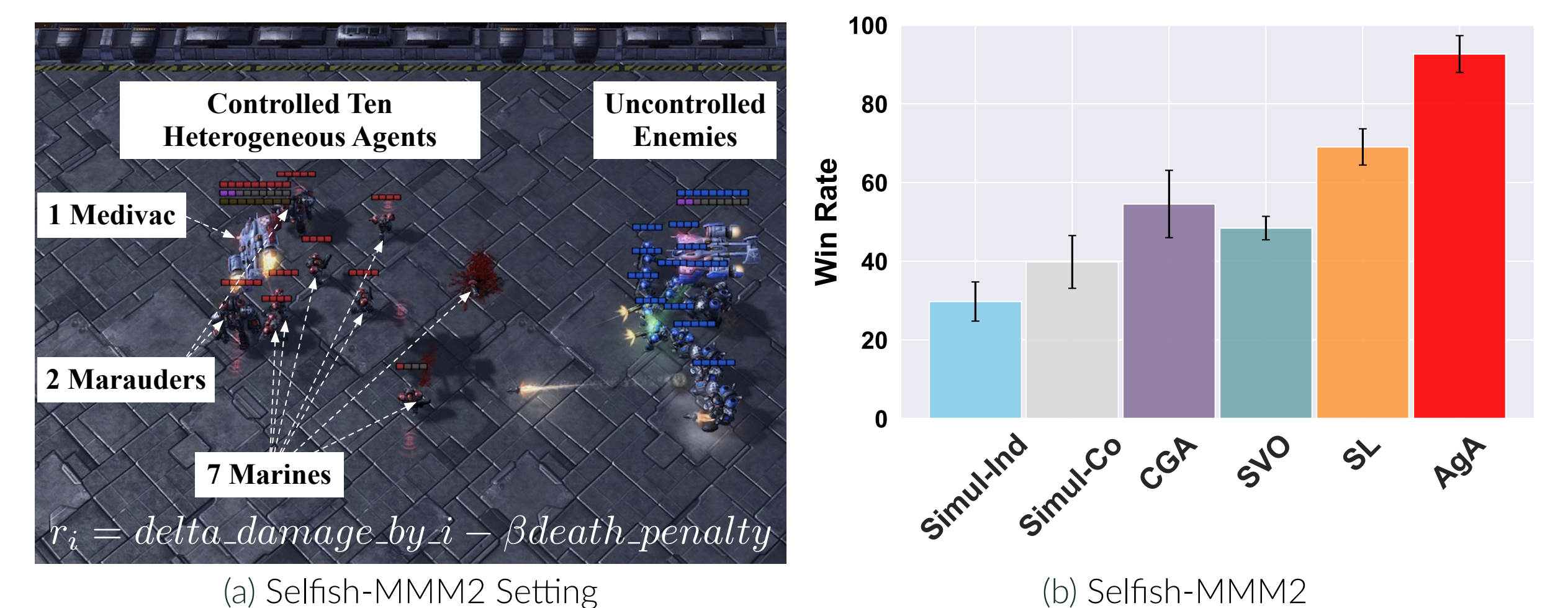
Cleanup and Harvest: common-used mixed-motive testbeds



(a) Cleanup

(b) Harvest

Selfish-MMM2: a large-scale mixed-motive testbed



(a) Selfish-MMM2 Setting

(b) Selfish-MMM2