

# Ranking Joint Policies in Dynamic Games using Evolutionary Dynamics

Anonymous Author(s)

Submission Id: «EasyChair submission id»

## ABSTRACT

Game-theoretic solution concepts, such as the *Nash equilibrium*, have been key to finding stable joint actions in multi-player games. However, it has been shown that the dynamics of agents' interactions, even in simple two-player games with few strategies, are incapable of reaching *Nash equilibria*, exhibiting complex and unpredictable behavior. Instead, evolutionary approaches can describe the long-term persistence of strategies and filter out transient ones, accounting for the long-term dynamics of agents' interactions. Our goal is to identify agents' joint strategies that result in stable behavior, being resistant to changes, while also accounting for agents' payoffs, in dynamic games. Towards this goal, and building on previous results, this paper proposes transforming dynamic games into their empirical forms by considering agents' strategies instead of agents' actions, and applying the evolutionary methodology  $\alpha$ -Rank to evaluate and rank strategy profiles according to their long-term dynamics. This methodology not only allows us to identify joint strategies that are strong through agents' long-term interactions, but also provides a descriptive, transparent framework regarding the high ranking of these strategies. Experiments report on agents that aim to collaboratively solve a stochastic version of the graph coloring problem. We consider different styles of play as strategies to define the empirical game, and train policies realizing these strategies, using the DQN algorithm. Then we run simulations to generate the payoff matrix required by  $\alpha$ -Rank to rank joint strategies.

## KEYWORDS

Evolutionary Dynamics, Empirical Games, Stochastic Games, Deep Reinforcement Learning, Ranking Joint Strategies

### ACM Reference Format:

Anonymous Author(s). 2025. Ranking Joint Policies in Dynamic Games using Evolutionary Dynamics. In *Proc. of the 24th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2025)*, Detroit, Michigan, USA, May 19 – 23, 2025, IFAAMAS, 10 pages.

## 1 INTRODUCTION

Game theory studies agents' strategies not only in terms of optimality of performance but also with regard to stability of agents' behavior. Game-theoretic solution concepts, particularly the *Nash equilibrium*, have played an important role in this research. However, solution concepts do not account for the long-term dynamics of agents' interactions, which are important in dynamic settings. In static games, where payoff matrices are known, studying solution

Table 1: Payoff matrix for the Rock-Paper-Scissors game.

	Rock	Paper	Scissors
Rock	0,0	-1,1	1,-1
Paper	1,-1	0,0	-1,1
Scissors	-1,1	1,-1	0,0

concepts is relatively straightforward. For example, consider the payoff matrix for the Rock-Paper-Scissors game in Table 1. The mixed strategy *Nash equilibrium* occurs when both players randomize their choices uniformly across Rock, Paper, and Scissors. In dynamic settings involving sequential decision making, one must account for the dynamics of agents' interactions over time. In these settings, we need to analyze agents' behavior in terms of their payoffs, identifying joint strategies that result into agents' stable behaviors. Evolutionary approaches have shown great potential towards this aim.

To study agents' behavior in multi-agent dynamic settings, researchers often train deep learning models to learn joint policies. These models, either in collaborative or competitive settings, are usually trained with the ultimate objective to result into Nash equilibria, aiming to agents' stability of behavior, where no agent has an incentive to deviate from their joint policy. In complex dynamic settings with long-term dynamics of agents' interactions, there is no guarantee of reaching that objective and there is no way to reveal the reasoning behind the agents' choice of a policy instead of another. Although proposals towards explainability and interpretability of models are important, these aim to provide either explanations for the policy as a whole (i.e. agent's style of play) or about individual decisions. In our case, we need a descriptive framework to account for transparency regarding the strength of agents' joint policies, accounting for long-term dynamics.

Our goal is to identify agents' joint policies that result in stable behavior, being resistant to changes, while also accounting for agents' payoffs, in dynamic games. This is motivated by the need to identify joint strategies of agents, whether human or software, that need to act as co-players in a common setting. To address this, we propose using a descriptive evolutionary framework that accounts for long-term agents' interactions. We conjecture that this approach helps agents select strong (i.e. non transient) policies when playing with or against other agents in dynamic settings where their actions affect future states and decisions. Identifying these joint policies and understanding their "superiority" is important, particularly in collaborative scenarios where agents must choose strategies while interacting with humans that use specific styles of problem-solving.

Towards this goal we propose exploiting multiple policy models, each realizing a distinct style of play (*strategy*), and then defining an empirical game for evaluating agents' joint performance when they play using various strategy profiles. This empirical game is

exploited by the  $\alpha$ -Rank evolutionary framework [5] to evaluate the evolutionary dynamics of agents' strategies over time, ultimately identifying which ones prevail in the long run.

Although this work builds on the  $\alpha$ -Rank framework, it contributes a perspective for evaluating individual agents' strategies in stochastic, sequential decision making settings, when they act with other agents following specific styles of play. In so doing, we do:

- Describe a concise methodology for evaluating and ranking agents' joined policies, accounting for their long-term interactions in dynamic settings, using the  $\alpha$ -Rank evolutionary framework.
- Demonstrate this methodology in multi-agent graph coloring dynamic games, defining multiple styles of play (i.e. strategies) per agent.
- Show how agents' choices of strategies—and the policies realizing them—can be transparently justified, by means of a descriptive framework.

## 2 BACKGROUND

In this section, we outline the key concepts necessary to follow the proposed approach.

### 2.1 Dynamic Games

Dynamic games describe agent interactions along the time dimension. As opposed to static games, where players execute single, one-shot actions, dynamic games involve a series of decisions made by each of the players at subsequent points in time. A key property of dynamic games is that the actions taken at any given moment influence the future states of the system and future decisions made. These temporal dependencies require players to consider the long-term consequences of their actions. A dynamic game can be represented as a tuple  $G = (S, K, A, T, P)$ , where  $S$  represents a finite set of states,  $K$  is the set of players, and  $A = (A^k \times A^{-k})$  is the set of joined actions, with  $A^k$  corresponding to the action set available to player  $k$ .  $A^{-k}$  denotes the action set available to players other than  $k$ . The transition matrix  $T$  describes how states evolve over time, determining the next state of the system based on the current state and the actions chosen by the players. Finally,  $P^k : S \times (A^k \times A^{-k}) \times S \rightarrow \mathbb{R}^K$  is the payoff function for player  $k$ , given the current joint state, the action chosen by player  $k$  and the actions of the other agents, and the resulting state.

In this work we focus on stochastic dynamic games, as introduced by L.S. Shapley in 1953 [7]. In stochastic games, the outcome of players' actions is influenced by probabilistic events, making future states of the game uncertain. These games are often referred to as Markov games [8]. Therefore, in stochastic games, the transition function  $T$  is defined as a probability distribution over next states. Specifically,  $T : S \times A \rightarrow \Delta(S)$ , where  $\Delta(S)$  is a probability distribution over the states, given a state and joint action. For example, in poker, while players' actions do influence the outcome, the next state of the game also depends on luck, such as drawing a strong hand like a flush or a weak hand like a pair of twos. In such games, players, when planning their actions, must account for both the actions of their opponents and the dynamics of the environment.

In dynamic games, players aim to decide on the course of their joint actions through time (joint policy) to maximize their accumulated discounted rewards over time:  $\sum_{s_{t+1} \in S} \gamma T(s_t, (a_t^k, a_t^{-k}), s_{t+1}) \cdot P^k(s_t, (a_t^k, a_t^{-k}), s_{t+1})$ , where  $\gamma \in (0, 1)$  is the discount factor that reflects the relative importance of future rewards compared to immediate rewards. Here,  $T$  represents the transition from state  $s_t$  to the state  $s_{t+1}$ , and  $P^k(s_t, (a_t^k, a_t^{-k}), s_{t+1})$  is the reward the player receives for choosing action  $a_t^k$ , given the actions  $a_t^{-k}$  of the other players, at state  $s_t$ , and resulting into state  $s_{t+1}$ .

### 2.2 Empirical Analysis and Empirical Games

Empirical Game Theory Analysis (EGTA) provides a framework that uses empirical methods to analyze player interactions within complex game environments [3]. These methods are used to define game components, such as payoff matrices, based on observed interactions, rather than relying on predefined rules. Simulation is one such method, where agents repeatedly play a game, and payoffs are collected based on the outcomes of these interactions. Other techniques include sampling, where a subset of the action space is explored to approximate the payoffs for a wider set of actions, and machine learning methods to identify players' behavior and estimate outcomes based on historical data [11]. Empirical techniques are applied in cases where the action space is too large and complex to define manually, making payoff matrices impossible to generate from simple rules and assumptions.

An empirical game, also referred to as a meta-game, is a *Normal Form Game* of the form  $G = (K, Str, P)$ , where  $K$ ,  $Str$  and  $P$  specify players, players' strategies, and payoffs, correspondingly. We define empirical games by abstracting the actions and defining the payoffs of players in an underlying dynamic game. The underlying game represents the actual setting where players interact. In the empirical game representation,  $K$  is the same as in the underlying game, i.e. the set of players engaged in strategic interactions. Strategies (i.e. styles of playing the game) in empirical games offer an action abstraction and can be derived by identifying distinct behaviors during game-play. The strategy space  $Str$  consists of distinct agents' styles of play.  $Str^k$  denotes the strategies of agent  $k$  and  $Str^{-k}$  the set of strategies of agents other than  $k$ . The set of strategy profiles, i.e. agents' joint strategies, is defined to be  $\mathcal{SP} = \{S_i | S_i = (str_i^1, str_i^2, \dots, str_i^K), \text{ where } str_i^k \in Str^k, \text{ and } i = 1, \dots, N \text{ the profile index}\}$ . The payoff matrix of an empirical game can be generated using empirical analysis techniques. Here, we focus on simulation, where agents engaged in the underlying game act according to policies adhering to specific strategies.

Subsequently, we use the terms *action* and *policy* when speaking about the underlying game, and the term *strategies* or *styles of play* when speaking about the empirical game.

The payoff function  $P$  of the empirical game is computed from simulations for each strategy profile as follows:

$$P^k(str^k, str^{-k}) = \frac{1}{N} \cdot \sum_{i=1}^N P_i^k(str^k, str^{-k})$$

where  $N$  is the number of simulation runs,  $str^k$  represents player  $k$ 's strategy,  $str^{-k}$  denotes the strategies of the other players, and

$p_i^k(str^k, str^{-k})$  (with an abuse of notation) represents the payoff player  $k$  receives in simulation run  $i$  when playing strategy  $str^k$  against the strategies of the other players. It must be noted that in contrast to dynamic games the payoff function does not take states as arguments, as the outcomes are determined by agents' joint strategies, i.e.  $P^k : (Str^k \times Str^{-k}) \rightarrow \mathbb{R}^K$  [5]. If we aggregate these expected payoffs into a matrix, we get the empirical payoff matrix whose dimensionality is  $\prod_{k=1}^K Str^k$ . Each entry represents the expected payoff for strategy  $str^k$  against strategy  $str^{-k}$ .

### 2.3 The $\alpha$ -Rank Method

Evolutionary dynamics studies how agents' interactions in multi-agent settings evolve over time. While single-agent systems have acquired a strong foundation over the years [1], multi-agent systems are more challenging to analyze.

Current literature indicates a growing interest in studying the evolutionary dynamics of multi-agent systems [1] [2] [6]. Although one might view evolutionary algorithms as mere tools for agents' hyper-parameter tuning, their contributions extend far beyond that. In the context of games, evolutionary algorithms are widely used to explore game-theoretic concepts. This area of study is also known as *Evolutionary Game Theory*. Building on work done in Evolutionary Game Theory,  $\alpha$ -Rank [5] introduces a novel game-theoretic approach to provide insights into the long-term dynamics of agents' interactions.

$\alpha$ -Rank is an evolutionary methodology designed to evaluate and rank agents' strategies in large-scale multi-agent interactions, using a new dynamic solution concept called *Markov-Conley chains* (MCCs). Given a K-player game,  $\alpha$ -Rank considers the empirical game with K player slots, called *populations*, where individual agents correspond to strategies, i.e. to styles of playing the underlying game.

In  $\alpha$ -Rank, populations of agents interact with each other through an evolutionary process following the dynamics of games. The rewards received from these interactions determine how well each strategy performs and, in turn, how often it is adopted by individuals in the populations. Strategies that perform well have a higher probability of being adopted and carried over to the next generation, while those performing poorly are less likely to be adopted. This process of competition and selection between populations leads to their evolution.

To facilitate evolution,  $\alpha$ -Rank uses the concept of mutation. Initially, populations are monomorphic, meaning all individuals within them choose the same strategy. During K-wise interactions, individuals have a small probability of mutating into different strategies or choosing to stick with their current one. The probability that the mutant will take over the population, defined to be the fixation probability function  $\rho$ , depends on the relative fitness of the mutant and the population being invaded. Fitness is a function that computes the expected reward an individual can receive when adopting a particular strategy, given the strategies of the other individuals. The stronger the fitness, the more likely it is for individuals to mutate, whereas the lower the fitness, the more likely it is for the mutant to go extinct. When the mutation rate is small, we can assume that the fitness for any agent  $k$  is  $f^k(str^k, str^{-k}) = P^k(str^k, str^{-k})$ , where  $P$  is the empirical game payoff.

Formally, the probability of a mutant strategy  $str'$  fixating in some population where individuals play strategy  $str$  is given by:

$$\rho_{str \rightarrow str'} = \frac{1 - e^{-\alpha \cdot \Delta f}}{1 - e^{-\alpha \cdot m \cdot \Delta f}} \quad (1)$$

assuming that  $\Delta f$  is non-zero.  $\Delta f = f^k(str', str^{-k}) - f^k(str, str^{-k})$  represents the difference in fitness between the mutant strategy  $str'$  and the resident strategy  $str$  in the focal population  $k$ , while the remaining  $K - 1$  populations are fixed in their monomorphic strategies  $str^{-k}$ . Parameter  $m$  is the population size and  $\alpha$  is the selection intensity. This adjusts the sensitivity of the system to fitness differences: with higher values of  $\alpha$ , even small differences in fitness lead to larger changes in  $\rho$ . The nominator measures the potential of the mutant to "invade" the resident population solely based on its fitness advantage. Note that, for example, as  $\Delta f$  approaches zero, the probability of the mutant's success decreases. The denominator, on the other hand, normalizes the fixation probability using the population size  $m$ , making it more challenging for a mutant to dominate in larger populations. When  $\Delta f$  is zero, the fixation probability comes down to  $1/m$ , indicating that the mutant strategy has the same probability of taking over as any other strategy in the population. We refer to this probability as the *neutral fixation probability*, denoted by  $\rho_m$ .

In the context of K-player games,  $\alpha$ -Rank creates a Markov transition matrix over strategy profiles. This is an  $|Str| \times |Str|$  matrix that defines the probability of moving from one strategy profile to another based on how likely each population is to change its strategy.

$$C_{str \rightarrow str'} = \begin{cases} \eta \cdot \rho_{str \rightarrow str'} & \text{if } str \neq str' \\ 1 - \sum_{str \neq str'} C_{str \rightarrow str'} & \text{otherwise} \end{cases} \quad (2)$$

Here,  $C$  is the strategy-transition matrix where each entry  $C_{str \rightarrow str'}$  represents the probability of transitioning from strategy  $str$  to strategy  $str'$ . The first part of the formula, calculates the probability of transitioning from one strategy to a different one, scaled to ensure that the sum of probabilities for all possible transitions from that strategy sums up to 1. The second part of the formula, computes the probability of staying with the same strategy,  $str$ , by excluding transitions to all other strategies.

This evolutionary process of competition and selection among players' strategies leads to a unique stationary probability distribution  $\pi$  of dimensionality  $|SP|$ , where the mass assigned to a strategy profile indicates how likely it is to resist being "invaded" by other strategies as the dynamics evolve. To evaluate and rank strategy profiles—which is the ultimate goal—the method calculates  $\pi$  over the game's Markov chain, using the strategy-transition matrix  $C$ . This distribution indicates how often the system is likely to remain in each profile over time, allowing us to identify the most dominant strategies that are expected to prevail in the long run. Formally,  $\pi$  can be computed from the following equation:

$$\pi C = \pi \Rightarrow \pi(C - \mathbb{I}) = 0 \quad (3)$$

where  $\mathbb{I}$  is the identity matrix. This means we are looking for a

probability vector  $\pi$  such that when multiplied by the transition matrix  $C$ , it remains unchanged. To solve for  $\pi$ , the augmented matrix from  $C - \mathbb{I}$  is constructed and a normalization condition to ensure that probabilities sum to 1 is imposed<sup>1</sup>. In this stationary distribution,  $\pi = (\pi_1, \pi_2, \dots, \pi_{|(\mathcal{SP})|})$ , each  $\pi_i$  represents the average time the system spends in strategy profile  $i$ .

### 3 PROBLEM STATEMENT

As already stated, we aim at identifying (human and software) agents' strong joint strategies, in terms of stability and joint performance, to solve problems in dynamic settings, accounting for agents' long-term dynamics of interactions. Stability implies non-transient strong strategies, persisting in time, as they fit better to the objective of the agents given the structure of the game and payoffs received. However, in dynamic games, we need to define the payoff matrix and exploit this to determine strategies stability. Even if we manage to estimate payoffs, the computation of solution concepts like the *Nash equilibrium* imposes a high computational cost in these settings, does not guarantee convergence, and fails to scale to large games. Beyond identifying stable joint strategies, it is important to transparently justify/describe what makes one joint strategy better than another. This requires more than just providing rankings of strategy profiles; it requires providing evidence for the rankings.

We could, therefore, consider our problem as follows: Given a dynamic game  $G$  with  $K$  players, our goal is to identify styles of playing  $G$ , and thus, the set of strategy profiles  $\mathcal{SP}$ , and rank these profiles based on how stable they are over time, considering long-term agents' interactions towards achieving their objectives. Specifically, we aim to define a ranking function  $\mathcal{R} : \mathcal{SP} \rightarrow \mathbb{R}$ , where  $\mathcal{R}(S_i) > \mathcal{R}(S_j)$  (resp.  $\mathcal{R}(S_i) \geq \mathcal{R}(S_j)$ ) indicates that the strategy profile  $S_i$  is strictly (resp. weakly) preferred over  $S_j$ , using a descriptive framework  $\mathcal{D} : \mathcal{SP} \times \mathcal{SP} \rightarrow \mathbb{R}$  that provides transparency on how rankings are decided.

It must be noted that empirical game strategies are realized by agents' policies adhering to these strategies in the underlying game. Thus, identifying stable joint strategies in the empirical game translates to identifying stable joint policies adhering to these strategies in the underlying dynamic game.

### 4 PROPOSED METHOD

To address the challenge of identifying stable joint policies in dynamic games, we propose an approach that combines concepts from *Empirical Game Theory* and *Evolutionary Dynamics*, using  $\alpha$ -Rank, providing transparency to rankings of agent's styles of play.

Given that the set of agents' policies in dynamic games can be infinitely large we focus on a subset of policies that adhere to concrete and well-defined styles of play. A way to identify styles of play is to observe how players behave in the underlying game or exploit demonstrations of game playing. For instance, human experts performing a task usually follow a distinct set of specific styles based on well-established practices, preferences and experience. Having

determined the game playing strategies, we can transform the dynamic game into its empirical form, defining the meta-game, as specified in Section 2.2: By (a) identifying empirical game strategies, and (b) training policies for agents to play the underlying game according to these strategies, (c) defining the empirical game payoff matrix, through simulations, exploiting the trained policies.

Having defined the meta-game, we need to define the function  $\mathcal{R}$ , which ranks joint strategies based on agents' long-term dynamics and objectives. In our approach, we propose using the evolutionary  $\alpha$ -Rank methodology to determine these rankings. The rankings are based on each strategy profile evolutionary success, which is reflected in the probability of that profile being selected over time. This probability is captured by the stationary distribution  $\pi$ , which  $\alpha$ -Rank computes in the limit of infinite ranking intensity  $\alpha$ . As demonstrated by [5], a large  $\alpha$  limit suffices. Therefore the long-term behavior is captured by the unique stationary distribution  $\pi$  under the large  $\alpha$  limit. As it is proved in [5], the Markov chain associated with a generalized multi-population model, coincides with the MCC solution concept. MCCs can be found efficiently in all games and can be identified by the sink strongly connected components of a response graph, whose vertices correspond to pure strategies' profiles and a directed edge from strategy profile  $S_i$  to a strategy profile  $S_j$  specifies that  $S_j$  is weakly a better response than  $S_i$  for player  $k$ .

To compute  $\pi$  over strategy profiles,  $\alpha$ -Rank requires the payoff matrix of the empirical game  $P$ . Along with the stationary distribution  $\pi$ ,  $\alpha$ -Rank also outputs the fixation probability function  $\rho_{S_i \rightarrow S_j}$ , where  $S_i, S_j \in \mathcal{SP}$ . One could abstractly illustrate  $\alpha$ -Rank as a function:

$$\alpha\text{-Rank}(P) \rightarrow (\pi, \rho) \quad (4)$$

While the stationary distribution  $\pi$  provides valuable insight into the long-term behavior of strategies, it alone does not help us fully understand how strategies transition between one another. The fixation probability function  $\rho$ , which measures the likelihood of transitioning from one strategy profile  $S_i$  to another  $S_j$ , fills this gap. Based on this, the descriptive framework  $\mathcal{D}$  can be adequately represented by  $\pi$  and  $\rho$ , which are constituents of the response graph, which provides a complete view of the empirical game dynamics.

Overall, building on the  $\alpha$ -Rank descriptive framework, the method proposed here for computing strategy profile rankings in dynamic games is as follows:

- (1) Identify players' styles of play.
- (2) Define the strategies of the empirical game based on those styles.
- (3) Train policies realizing the defined strategies.
- (4) Run game simulations to create the empirical payoff matrix  $P$ .
- (5) Apply  $\alpha$ -Rank to define  $\mathcal{R}$  and  $\mathcal{D}$ :
  - (a) Calculate the Markov transition matrix  $C$ .
  - (b) Find the unique stationary distribution  $\pi$ .
  - (c) Rank joint strategies by ordering the masses of  $\pi$ .
  - (d) Describe the rankings through the response graph.
  - (e) Study the effect of different  $\alpha$  values on  $\pi$ .

<sup>1</sup>The system  $\pi(C - \mathbb{I}) = 0$  by itself does not have a unique solution, as there are infinitely many vectors  $\pi$  that satisfy it. To get a unique solution  $\pi = (\pi_1, \pi_2, \dots, \pi_{|(\mathcal{SP})|})$ , it must hold that  $\sum_i \pi_i = 1$ .

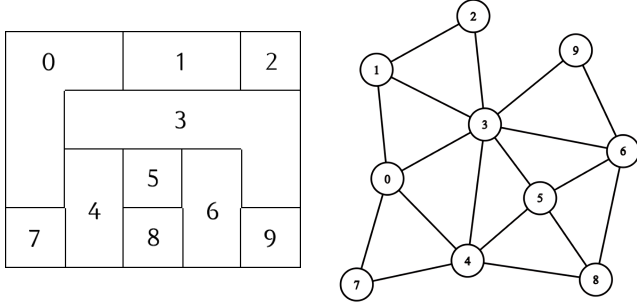
## 5 EXPERIMENTS AND RESULTS

### 5.1 The Graph Coloring Problem

The *Graph Coloring Problem* (GCP) is one of the most well-known problems in graph theory. It involves assigning colors to vertices in a graph such that no two adjacent vertices share the same color, and using the minimum number of colors, also known as the chromatic number [10].

In this study, we shift our focus from finding the chromatic number across graph configurations to solving the multi-agent problem of assigning colors to vertices of a dynamic graph with respect to the constraints: In doing so, we define the graph coloring problem as a dynamic game that allows us to study the evolutionary dynamics in multi-agent interactions.

Through this problem, we aim to demonstrate how we can gain insights into individuals’ joint strategies, i.e. into the effectiveness of playing the dynamic game when individual styles of play are combined.



**Figure 1: A snapshot of the game environment in grid and graph forms.**

We consider the underlying graph coloring dynamic game to be a two-player game executed in rounds. The graph corresponds to a grid comprising blocks of cells: A block comprises one or more merged cells. Each vertex of the graph corresponds to a block, and the adjacency relation between blocks specifies the edges in the graph. At the beginning of the game, the grid is initialized with a random number of rows and columns ( $n \times m$ ). In our experimental setup we assume a  $4 \times 5$  grid. The environment is initialized by randomly combining cells to create the blocks. The resulting configuration remains the same throughout the entire game. A snapshot of such a configuration with 10 blocks is the one shown in Figure 1, together with the corresponding graph. Merging cells is important as it allows for complex neighboring relationships to be defined, expanding beyond the standard constraints between adjacent blocks. Blocks are either (a) colored by the agents, (b) white (free to be colored) or (c) hidden (their colors cannot be observed and they can not be re-colored by the agents). Let  $B$  the set of blocks corresponding to graph vertices, and  $CR$  the set of possible colors that an agent can use for coloring blocks in  $B$ . The game unfolds over multiple rounds in which agents simultaneously choose their actions. At the beginning of each round, the environment reveals the color of some of the hidden blocks, if any. The number of blocks that get un-hidden is random, which implies that the state of the

Preference Dimension	Value
Color Tone	warm (W) vs. cool (C)
Block Coloring Difficulty	small (L) vs. large (A)
Coloring Approach	minimalistic (M) vs. extravagant (E)

**Table 2: Dimensions specifying agents’ strategies**

graph is influenced not only by the agents’ actions, but also by the environment. We therefore consider the game to be stochastic. As soon as all blocks in  $B$  are uncovered and colored, the game ends.

The set of agents’ actions  $A$  is defined to be the Cartesian product of the set of blocks  $B$  and the set of the available colors  $CR$ , denoted as:

$$A = B \times CR = \{(b, c) \mid b \in B, c \in CR\} \quad (5)$$

To specify states, let  $CR^*$  include the elements of  $CR$ , and two additional elements representing hidden and white blocks:  $CR^* = CR \cup \{\text{hidden}, \text{white}\}$ . A state  $s$  is as follows:

$$s = \{(b_i, c_i), i = 1, \dots, |B|\}, \\ s.t. \forall b \in B, \exists \text{ a unique } c \in CR^*, \text{ with } (b, c) \in s$$

Regarding the reward function, it is a sum of gains, penalties, sanctions, delays and adopted preferences. Given that actions are represented as vectors of shape  $(b, c) \in B \times CR$ , an agent receives a gain point (+1) for each neighbor of  $b$  that has a different color than the chosen color  $c$ . On the contrary, an agent receives a penalty point (-2) for each neighbor that shares the same color  $c$ . Sanction is a big negative reward (-10) that an agent receives when it attempts to color a hidden block or a block that has already been colored. Delay (-1) is a small negative reward that both agents receive when they try to color the same block  $b$ , causing a brief pause in the game to determine which agent will eventually color  $b$ . Last but not least, there is the preference-adoption reward, which agents receive regardless of whether their action is good, bad, or forbidden. This reward helps agents to be trained so as to adhere to specific preferences, or what we call *styles of play*. We will elaborate shortly on these in the following section.

### 5.2 Defining the Empirical Game

Transforming the underlying dynamic graph-coloring game into its empirical form involves two key steps: (1) identifying agents’ strategies and (2) constructing the empirical game payoff matrix.

**5.2.1 Agents’ Strategies.** Agents’ strategies define distinct styles of play, usually revealed by preferences in playing the game. In our experiments, we specify different styles across three main dimensions: color tone (preference for which colors to use), block difficulty (preference for the types of blocks to choose), and coloring approach (preference for the number of colors to use), as shown in Table 2. Through the combination of preferences in each of these dimensions, a style can range from complete indifference, where none of the dimensions hold any influence (denoted by “I”), to specific preferences in all dimensions.

Policies corresponding to specific strategies are represented using convolutional neural networks. To train these policies, we assign specific values in the three dimensions of the game’s *preference* reward. These values, range from -1 to 1, where 1 indicates a strong preference for a particular dimension. For example, a value of 0.7 for warm colors indicates a relatively high preference for warm tones. In our experimental setting, we define 11 distinct styles: I, C, W, E, M, L, A, AE, CA, LE and WL, given “I” and combinations of preference values specified in Table 2.

Assuming no inherent bias among the players of the empirical game, we allow populations to sample from the same list of strategies.

**5.2.2 Training the agents.** All policy models share the same underlying architecture and training setup. Although hyperparameter tuning is typically recommended, it doesn’t make much difference in this case, as these models are relatively easy to optimize when trained in small settings. Regarding the convolutional neural network architecture, it consists of four convolutional layers, each defined with a kernel size of 3, stride of 1, and padding of 1, meant to extract spatial features from the input. The input tensor has dimensions  $10 \times 12$ , where  $|B| = 10$  represents the number of blocks in the state and  $|CR^*| = 12$  represents the number of possible colors a block can have. Each block is encoded using one-hot encoding, meaning that each color is represented as a binary vector of length 12. The output is then flattened and passed through two fully connected layers, which process the data to produce the final output, as shown in Figure 2.

Policy models are trained individually (with no co-players) in the underlying game using the deep Q-learning reinforcement learning algorithm specified in Algorithm 1. We set  $\gamma$  to 0.7. To optimize the model parameters, we use the smooth L1 loss function with  $\beta=1.0$  and the Adam optimizer with a learning rate of  $5e-4$  and weight decay of  $1e-5$  to prevent over-fitting. To further enhance the learning process, we incorporate experience replay, with a memory that stores up to 10 million experiences [4]. A target network alongside the main policy network, is being used according to the Double-DQN approach [9]. To update the target network we apply a soft update with a factor  $\tau=5e-3$ . This gradually brings the target network closer to the policy network, balancing learning speed and stability. With a batch size of 64, we train the models for 10000 episodes.

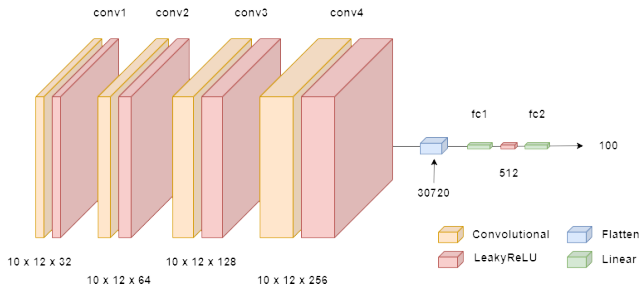


Figure 2: Convolution Policy Network Architecture

Algorithm 1: Double Deep Q-Learning with Experience Replay

```

1:  $Q_\theta, Q_{\theta'} \leftarrow Q_\theta, M$   $\triangleright$  Initialize policy/target nets & memory
2: for episode do
3:    $s \leftarrow s_0$ 
4:   for step do
5:      $a \leftarrow \operatorname{argmax}_a Q_\theta(s)$   $\triangleright$  Select  $\epsilon$ -greedy action
6:      $(s, a, r, s') \in M$   $\triangleright$  Store experience
7:     if  $|M| > \text{batch size}$  then
8:       for each  $(s, a, r, s')$  in  $M$  do  $\triangleright$  Sample memory
9:          $y \leftarrow r + \gamma \max_{a'} Q_{\theta'}(s')$ 
10:         $L \leftarrow \text{Loss}(Q_\theta(s), y)$ 
11:         $\theta \leftarrow \theta - \alpha \nabla_\theta L$ 
12:      end for
13:    end if
14:     $Q_{\theta'} \leftarrow \tau Q_\theta + (1 - \tau) Q_{\theta'}$   $\triangleright$  Soft update
15:     $s \leftarrow s'$ 
16:  end for
17: end for

```

**5.2.3 Empirical Game Payoff Matrix.** We generate the empirical payoff matrix by simulating each strategy profile over multiple games. These payoffs represent how well different styles of play perform jointly, according to the game’s rules.

The values in the payoff matrix are computed in terms of the delay and the quality of the solution according to the game’s constraints (gain, penalty, and sanction), excluding preferences. This ensures a common ground for distinct strategies, evaluating solutions solely based on the game’s rules. For each pair of strategies, we simulate the game over 5,000 repeats and calculate the average payoff for each strategy. These values are then organized into the payoff matrix, which is provided in Table 4 (Appendix A). From this matrix, we observe that (L, WL) and its symmetric counterpart (WL, L) both with payoffs of (3.15, 3.21) and (3.21, 3.15) respectively, are the only Nash equilibria. It is important to note here that these equilibria prescribe agents’ strategies given that they do play the game with rational co-players, but they do not capture the overall dynamics of the game, considering the long-term effects of agents’ interactions.

### 5.3 Evaluation and Ranking

Given the payoff matrix derived from the empirical analysis, we apply the  $\alpha$ -Rank method to evaluate the performance of strategy profiles over time in terms of the MCC solution concept. Specifically, we ran the method 1000 times, using values of  $\alpha$  within the range  $[0.1, 10]$  with step=0.01, while assuming populations of size  $m = 100$ . We provide as input the strategies defined in Section 5.2.1 and the empirical game payoff matrix. We focus on the rankings of the top 6 strategy profiles, to identify the stronger ones across different values of  $\alpha$ .

As we observe from the rankings in Table 3, the strategy profile that prevails in the long run is (WL, CA); this is the primary component of the MCC. Although the table was derived using an  $\alpha$  value of 2, the rankings remain consistent even when  $\alpha$  is set to 10. We choose  $\alpha = 2$  over  $\alpha = 10$ , to display the rankings of lower-performing strategy profiles, which would otherwise drop to



Agent	Rank	Score
(WL, CA)	1	0.42
(W, CA)	2	0.13
(M, CA)	3	0.12
(CA, M)	4	0.08
(CA, W)	5	0.08
(CA, LE)	6	0.01

**Table 3: Rankings for  $\alpha = 2$**

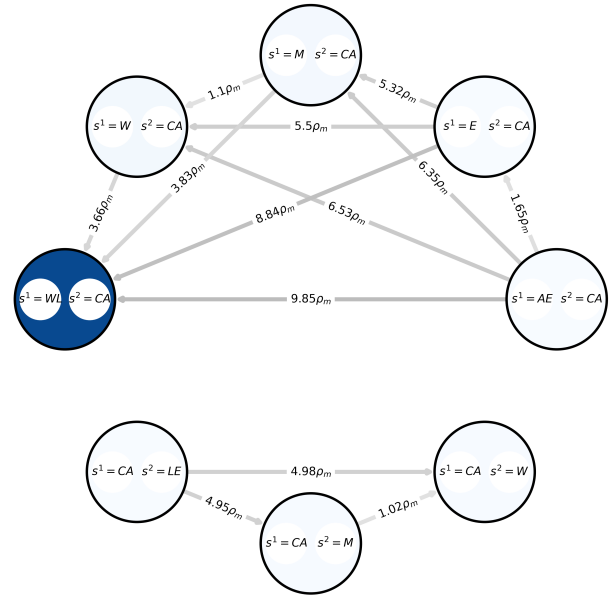
zero. First, it is worth mentioning that the Nash equilibria (L, WL) and (WL, L) don't appear among the top-ranked strategy profiles. This is because MCC components are defined based on how well strategies perform when interacting with other strategies, based on long-term agents interactions. The individual strategies within the Nash equilibrium profile, either WL or L, may not result in favorable interactions with other strategies. As a result, the profile (WL, L) is ranked lower than others.

To further support our observations regarding the misalignment between the two solution concepts, let's examine why (CA, WL) is part of the MCCs, while (L, WL), the Nash equilibrium, is not. A closer look at the payoff matrix in Table 4 reveals that L appears to be the worst-performing strategy for the row player, with an average payoff of 3.13. In this case, being in the Nash equilibrium means the player is stuck with a strategy that gives low rewards, making it the best among other options, rather than a strong choice. If it happens to play this strategy, it would expect its rational opponent to play WL. Strategy CA on the other hand, is the best-performing strategy for the row player, with an average payoff of 3.18. Combined with WL, which is the best performing strategy for the column player, with an average payoff of 3.18, they make profile (CA, WL) becomes the top ranked strategy profile in the ranking Table 3.

Rankings within the MCC are also very intuitive. For example, strategies that prefer different color tones, such as (WL, CA) or (W, CA), tend to result into fewer conflicts since, they naturally avoid selecting the same colors. Similarly, strategies that prefer different blocks based on their difficulty, such as (WL, CA) or (CA, LE), tend to provide solutions with minimal delay, as they naturally avoid coloring the same blocks. Notably, profiles with mixed preferences across these dimensions demonstrate the most promising performance, which explains why (WL, CA), as such a profile, is a key component of the MCC. However, not all profile rankings can be easily explained through the game's rules alone; the expected influence of certain strategies on the quality of the solutions remains ambiguous. For example, profiles with strategies like M and E are more difficult to analyze.

The response graph provides a visualization to interpret the  $\alpha$ -Rank results. This graph illustrates the MCC, using the strategy profiles' masses from the stationary distribution,  $\pi$ , along with the fixation probability function  $\rho$  provided by  $\alpha$ -Rank. Figure 5d shows the response graph for  $\alpha = 6.4$ . We consider it to be part of the descriptive framework  $\mathcal{D}$ , as it offers insights into how rankings were derived. Additional graphs for  $\alpha = 0.4, 1.3$ , and  $1.9$  are available in Figure 5 (Appendix B).

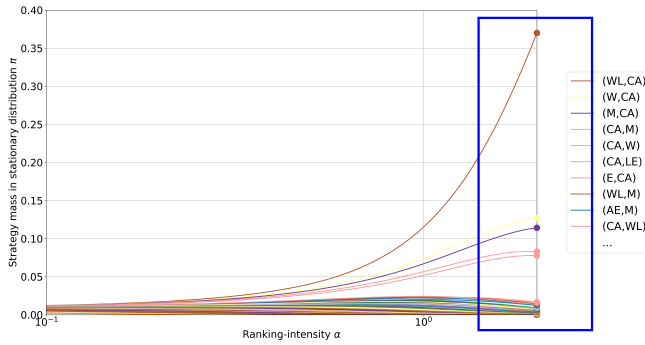
The response graph describes the overall dynamics of the strategy profiles in the empirical game. One prominent feature is the primary component of the MCC, specifically the profile (WL, CA). This profile, indicated by a dark blue color, has multiple graph edges leading to it, while none from it, indicating that strategies in this profile are non-transient. This is further supported by the large fixation probabilities along the edges. A particularly prominent example is the cluster (CA, LE)-(CA, M)-(CA, W), which consists of three strongly connected profiles, indicating that once a player adopts one of these profiles, they will likely remain within their cluster. These components reflect stable regions in the game's strategy dynamics, where transitions between profiles become locked into a cycle.



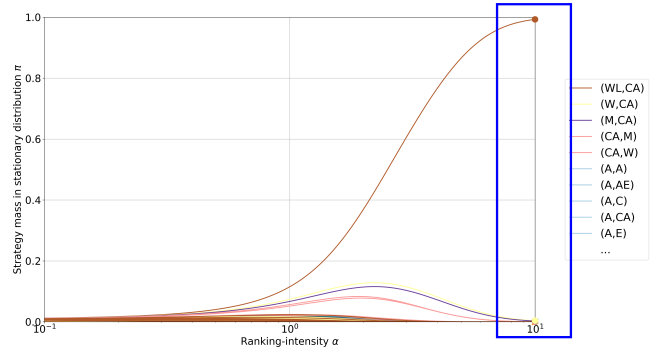
**Figure 3: Response graph for  $\alpha = 6.4$ .**

To further investigate the effect of  $\alpha$  on profile dominance, we plotted the stationary distribution  $\pi$  across all  $\alpha$  values used in the experiments, for the top-performing strategy profiles (see Figure 4). This visualization -also part of  $\mathcal{D}$ - helps us understand how the stationary distribution changes as the selection intensity increases. The x-axis represents the different  $\alpha$  values, ranging from 0.1 to 3 in Figure 4a, and from 0.1 to 10 in Figure 4b, while the y-axis in both figures shows the mass of each strategy profile in the stationary distribution  $\pi$ . As  $\alpha$  increases, the distribution converges, indicating that the selection process stabilizes. The final mass distributions are highlighted in boxed regions. The legend on the right side of the plot displays the top-performing joint strategies, with the stronger ones appearing at the top.

We plot two such graphs to observe how the mass of strategy profiles is distributed in the MCCs across different  $\alpha$  values. In the stationary distribution resulting from a bigger  $\alpha$ , the dominant strategy profile (WL, CA) in the MCC achieves a mass of 1, with all other profiles dropping to 0. This is clearly illustrated in the second



(a) Mass across  $\alpha \in [0.1, 3]$



(b) Mass across  $\alpha \in [0.1, 10]$

Figure 4: Effect of ranking intensity  $\alpha$  on strategy profile mass in the stationary distribution  $\pi$ .

plot (see Figure 4b). However, regarding the mass distribution for a smaller range of  $\alpha$ , depicted in the first plot, the game has not yet converged to the final MCC.

## 6 CONCLUSIONS

In this study, we developed a methodology for identifying strong joint-strategies in dynamic multi-agent games, accounting for stability and performance, using the  $\alpha$ -Rank evolutionary algorithm. The methodology is applied on a stochastic version of the *Graph Coloring Problem*, in which players work together to color a graph while ensuring that neighboring vertices are assigned different colors. According to the methodology, first we transformed the game into its empirical form, by defining strategies (styles of play). We then designed and trained Deep Q-Learning policy models that realize those styles of play in the underlying game, and run simulations to generate the empirical payoff matrix.  $\alpha$ -Rank, applied to this matrix, results into a unique stationary distribution over strategy profiles that defines the empirical game's MCC. The  $\alpha$ -Rank not only helped us identify stable strategy profiles resistant to changes but also provided a descriptive framework for understanding why certain profiles prevail in the long run, based on the underlying dynamics of the game. Through this approach, we successfully described a concise methodology for evaluating and ranking agents' joint policies, considering their long-term interactions in dynamic settings, while also explaining how strategy profiles are defined within the MCC.

Future work involves (a) applying the methodology in more complex and large-scale settings, accounting for strategy profiles of multiple stakeholders that may collaborate and/or compete, (b) using machine learning methods to identify different styles of play

from demonstrations and specifying the empirical game, (c) exploring advanced models able to adapt their strategies based on observed behaviors based on the behavior of co-players, and (d) applying the methodology into real-world settings where agents need to align with human preferences in dynamic settings.

## REFERENCES

- [1] Daan Bloembergen, Karl Tuyls, Daniel Hennes, and Michael Kaisers. 2015. Evolutionary dynamics of multi-agent learning: a survey. *J. Artif. Int. Res.* 53, 1 (May 2015), 659–697.
- [2] Omid E. David, H. Jaap van den Herik, Moshe Koppel, and Nathan S. Netanyahu. 2014. Genetic Algorithms for Evolving Computer Chess Programs. *IEEE Transactions on Evolutionary Computation* 18, 5 (Oct. 2014), 779–789. <https://doi.org/10.1109/tevc.2013.2285111>
- [3] Michael Levét. 2016. Game Theory : Normal Form Games. <https://api.semanticscholar.org/CorpusID:131771375>
- [4] Long-Ji Lin. 1992. Self-Improving Reactive Agents Based on Reinforcement Learning, Planning and Teaching. *Mach. Learn.* 8, 3–4 (may 1992), 293–321. <https://doi.org/10.1007/BF00992699>
- [5] Shayegan Omidshafiei, Christos Papadimitriou, Georgios Piliouras, Karl Tuyls, Mark Rowland, Jean-Baptiste Lespiau, Wojciech M. Czarnecki, Marc Lanctot, Julien Perolat, and Remi Munos. 2019.  $\alpha$ -Rank: Multi-Agent Evaluation by Evolution. arXiv:1903.01373 [cs.MA]
- [6] Sheryl Paul and Jyotirmoy V. Deshmukh. 2022. Multi Agent Path Finding using Evolutionary Game Theory. arXiv:2212.02010 [cs.MA] <https://arxiv.org/abs/2212.02010>
- [7] Lloyd S. Shapley. 1953. Stochastic Games\*. *Proceedings of the National Academy of Sciences* 39 (1953), 1095 – 1100. <https://api.semanticscholar.org/CorpusID:263414073>
- [8] Yoav Shoham and Kevin Leyton-Brown. 2008. *Multiagent Systems: Algorithmic, Game-Theoretic, and Logical Foundations*. Cambridge University Press.
- [9] Hado van Hasselt, Arthur Guez, and David Silver. 2015. Deep Reinforcement Learning with Double Q-learning. arXiv:1509.06461 [cs.LG]
- [10] George Watkins, Giovanni Montana, and Juergen Branke. 2023. Generating a Graph Colouring Heuristic with Deep Q-Learning and Graph Neural Networks. arXiv:2304.04051 [cs.LG]
- [11] Michael P. Wellman, Karl Tuyls, and Amy Greenwald. 2024. Empirical Game-Theoretic Analysis: A Survey. arXiv:2403.04018 [cs.GT] <https://arxiv.org/abs/2403.04018>



## A EMPIRICAL PAYOFF MATRIX

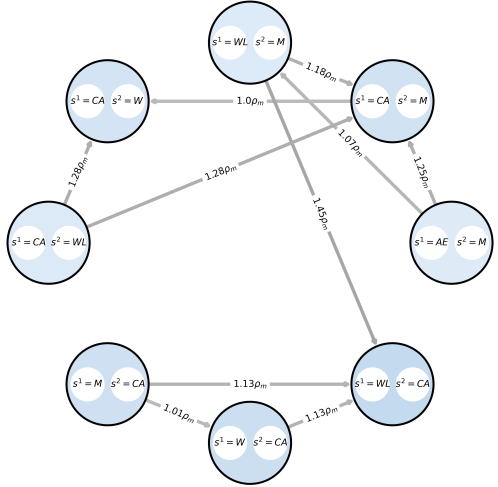
This is the empirical payoff matrix derived from simulations of the *Graph Coloring Game* using policies trained to adhere to specific styles of play. Each entry in the matrix represents the payoffs of strategies in the corresponding profile, with the first value indicating the payoff of the row strategy and the second value of the column player. The Nash equilibria are highlighted in bold, while nine of the top-ranked strategy profiles in the MCC are shaded in gray.

	A	AE	C	CA	E	I	L	LE	M	W	WL
A	(3.12, 3.11)	(3.15, 3.16)	(3.17, 3.17)	(3.14, 3.17)	(3.16, 3.17)	(3.16, 3.15)	(3.22, 3.13)	(3.19, 3.16)	(3.15, 3.18)	(3.16, 3.17)	(3.21, 3.18)
AE	(3.17, 3.17)	(3.11, 3.11)	(3.18, 3.17)	(3.15, 3.17)	(3.17, 3.16)	(3.19, 3.16)	(3.23, 3.12)	(3.19, 3.16)	(3.15, 3.18)	(3.17, 3.17)	(3.20, 3.16)
C	(3.17, 3.16)	(3.16, 3.17)	(3.10, 3.10)	(3.14, 3.17)	(3.15, 3.15)	(3.18, 3.15)	(3.22, 3.12)	(3.17, 3.14)	(3.14, 3.17)	(3.17, 3.16)	(3.20, 3.17)
CA	(3.17, 3.15)	(3.17, 3.15)	(3.17, 3.14)	(3.11, 3.11)	(3.18, 3.15)	(3.18, 3.14)	(3.24, 3.13)	(3.21, 3.16)	(3.16, 3.16)	(3.19, 3.16)	(3.22, 3.15)
E	(3.15, 3.16)	(3.16, 3.16)	(3.15, 3.16)	(3.15, 3.17)	(3.10, 3.10)	(3.18, 3.16)	(3.22, 3.12)	(3.19, 3.14)	(3.15, 3.17)	(3.16, 3.17)	(3.19, 3.17)
I	(3.14, 3.16)	(3.16, 3.18)	(3.16, 3.18)	(3.15, 3.19)	(3.16, 3.17)	(3.12, 3.12)	(3.22, 3.14)	(3.18, 3.16)	(3.14, 3.19)	(3.16, 3.18)	(3.19, 3.18)
L	(3.14, 3.22)	(3.11, 3.22)	(3.12, 3.22)	(3.13, 3.23)	(3.12, 3.22)	(3.13, 3.22)	(3.12, 3.12)	(3.14, 3.20)	(3.11, 3.21)	(3.14, 3.23)	<b>(3.15, 3.21)</b>
LE	(3.15, 3.19)	(3.14, 3.18)	(3.14, 3.18)	(3.15, 3.21)	(3.15, 3.19)	(3.16, 3.17)	(3.20, 3.14)	(3.11, 3.11)	(3.14, 3.22)	(3.15, 3.18)	(3.18, 3.19)
M	(3.17, 3.14)	(3.17, 3.15)	(3.17, 3.15)	(3.16, 3.17)	(3.16, 3.14)	(3.18, 3.14)	(3.23, 3.11)	(3.20, 3.14)	(3.06, 3.08)	(3.18, 3.15)	(3.20, 3.16)
W	(3.17, 3.17)	(3.17, 3.18)	(3.16, 3.18)	(3.16, 3.20)	(3.17, 3.17)	(3.18, 3.16)	(3.21, 3.13)	(3.18, 3.15)	(3.15, 3.18)	(3.08, 3.09)	(3.19, 3.15)
WL	(3.17, 3.20)	(3.17, 3.19)	(3.17, 3.19)	(3.17, 3.22)	(3.17, 3.19)	(3.18, 3.19)	<b>(3.21, 3.15)</b>	(3.19, 3.17)	(3.16, 3.20)	(3.16, 3.19)	(3.13, 3.13)

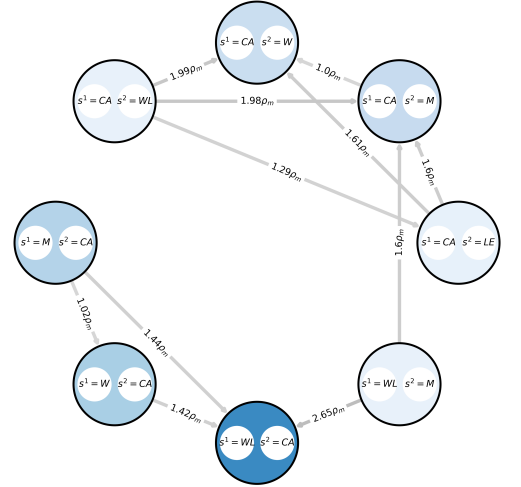
Table 4: Empirical Payoff Matrix for the Graph Coloring Game

## B RESPONSE GRAPH

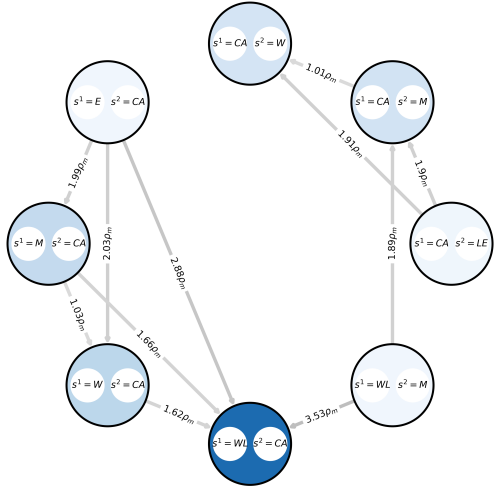
These are four response graphs illustrating the dynamics of strategy profiles in the empirical *Graph Coloring Game* for different  $\alpha$  values. Each node in the graph represents a unique strategy profile in the MCC, while the edges indicate transitions between them. The values on the edges show the fixation probabilities normalized by the neutral fixation probability, denoted as  $\rho_m$ . The nodes and edges are color-coded. Darker blue nodes represent more strong joint profiles, while lighter blue nodes represent transient ones. Similarly, bold arrows suggest a strong advantage in shifting between the nodes, whereas faint ones suggest less of an advantage.



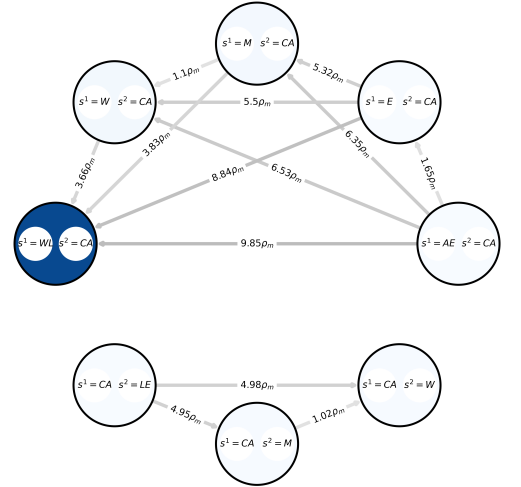
(a)  $\alpha = 0.4$



(b)  $\alpha = 1.3$



(c)  $\alpha = 1.9$



(d)  $\alpha = 6.4$

Figure 5: Response graphs of strategy profiles' dynamics.