

WHO TO IMITATE: IMITATING DESIRED BEHAVIOR FROM DIVERSE MULTI-AGENT DATASETS

Anonymous authors

Paper under double-blind review

ABSTRACT

AI agents are commonly trained with large datasets of demonstrations of human behavior. However, not all behaviors are equally safe or desirable. Desired characteristics for an AI agent can be expressed by assigning desirability scores, which we assume are assigned to collective trajectories, but not to individual behaviors. For example, in a dataset of vehicle interactions, these scores might relate to the number of incidents that occurred. We first assess the effect of each individual agent’s behavior on the collective desirability score, e.g., assessing how likely an agent is to cause incidents. This allows us to afterward only imitate agents with desired behavior, e.g., only imitating agents that are unlikely to cause incidents. To enable this, we propose the concept of an agent’s *Exchange Value*, which quantifies an individual agent’s contribution to the collective desirability score. This is expressed as the expected change in desirability score when substituting the agent for a randomly selected agent. We propose additional methods for estimating Exchange Values from real-world datasets, enabling us to learn aligned imitation policies that outperform relevant baselines.

1 INTRODUCTION

Imitating human behaviors from large datasets is a promising technique for achieving human-AI and AI-AI interactions in complex environments (Carroll et al., 2019; , FAIR; He et al., 2023; Shih et al., 2022). However, such large datasets can contain undesirable human behaviors, making direct imitation problematic. Rather than imitating all behaviors, it may be preferable to ensure that AI agents imitate behaviors that align with predefined desirable characteristics. In this work, we assume that desirable characteristics are quantified as desirability scores given for each trajectory in the dataset. This is commonly the case when the evaluation of the desirability of individual actions is impractical or too expensive (Stiennon et al., 2020). For complex datasets that involve multiple interacting agents, assigning desirability scores to collective trajectories – but not to individual behavior – may be the only viable option. For instance, in a football match, while the final score directly gauges team performance, determining individual player contributions is more difficult.

We develop an imitation learning method for multi-agent datasets that ensures alignment with desirable characteristics – expressed through a Desired Value Function (DVF¹) that assigns a score to each *collective* trajectory. This scenario is applicable to several areas that involve learning behavior from data of human groups. One example is a dataset of vehicle interactions and desirability scores which indicate the number of occurred incidents in a collective trajectory and the aim to imitate only behavior that is unlikely to result in incidents (e.g. aiming to imitate driving with foresight). Another example is a dataset of a multi-player online game and desirability scores reflecting players’ average enjoyment in each round and the goal to imitate only behavior that creates a positive experience.

Assessing the desirability of an individual agent’s behavior involves gauging its impact on the *collective* desirability score. For instance, it requires evaluating whether an agent’s behavior increases the likelihood of accidents while driving or decreases the enjoyment of other players in a game. This is termed the *credit assignment problem* (Shapley, 1953), akin to fairly dividing the value produced

¹The DVF itself is not sufficient to describe desired behavior completely, as it possibly only covers a subset of behavior, e.g., safety-relevant aspects. It is complementary to the more complex and nuanced behaviors that are obtained by imitating human demonstrations, providing guardrails or additional guidance.

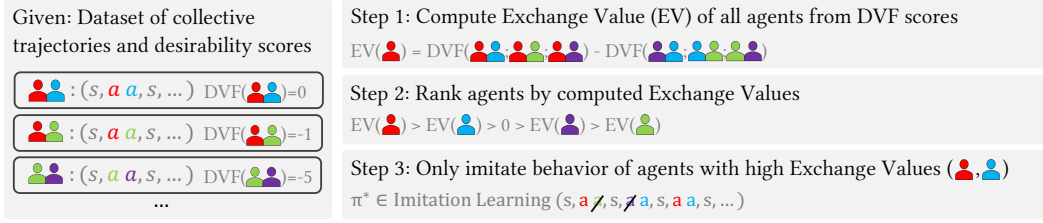


Figure 1: We are given a dataset composed of multi-agent trajectories generated by many individual agents, e.g., a dataset of cars driving in urban environments. In addition, the Desired Value Function (DVF) indicates the desirability score of a collective trajectory, e.g., the number of incidents that occurred. We first compute the Exchange Value (EV) of each agent, where a positive EV indicates that an agent increases the desirability score (e.g. an agent driving safely). We reformulate imitation learning to take into account the computed EVs, and achieve an imitation policy that is aligned with the DVF (e.g. only imitating the behavior of safe drivers).

by a group of players among the players themselves. The credit assignment problem proves complex in real-world scenarios due to three main factors (see Figure 2 for details): First, many scenarios only permit specific group sizes. This makes Shapley Values (Shapley, 1953) – a concept commonly used in Economics for credit assignment – inapplicable, as it relies on the comparisons of groups of different sizes. Second, real-world datasets for large groups are in practice always incomplete, i.e. do not contain trajectories for all (combinatorially many) possible groups of agents. Third, datasets of human interactions may be *fully anonymized* by assigning one-time-use IDs. In this case, if an agent is present in two trajectories, it will appear in the dataset as if it is *two different agents*, making the credit assignment problem degenerate. This requires incorporating behavior information.

To address these challenges we propose Exchange Values (EVs), akin to Shapley Values, which quantify an agent’s contribution as the expected change in desirability when substituting the agent randomly. EVs are applicable to scenarios with fixed group sizes, making them more versatile. We introduce EV-Clustering that estimates EVs from incomplete datasets by maximizing inter-cluster variance. We show a theoretical connection to clustering by *unobserved* individual contributions and adapt this method to fully-anonymized datasets, by considering low-level behavioral cues.

By incorporating agents’ estimated EVs, we introduce Exchange Value based Behavior Cloning (EV2BC), which imitates large datasets by only imitating the behavior of agents with EVs higher than a tuneable threshold (see Figure 1). This approach allows learning from interactions with agents with all behaviors, without necessarily imitating them, which is not possible when simply excluding all trajectories with a low collective desirability score. Our work makes the following contributions:

- We introduce *Exchange Values* (Def. 4.1) to compute an agent’s individual contribution to a collective value function and show their relation to Shapley Values.
- We propose *EV-Clustering* (Def. 4.4) to estimate contributions from incomplete datasets and show a theoretical connection to clustering agents by their unobserved individual contributions.
- We empirically demonstrate how EVs can be estimated from fully-anonymized data and employ EV2BC (Def. 4.5) to learn policies aligned with the DVF, outperforming relevant baselines.

2 RELATED WORK

Most previous work on aligning AI agents’ policies with desired value functions either relies on simple hand-crafted rules (Xu et al., 2020; , FAIR), which do not scale to complex environments, or performs postprocessing of imitation policies with fine-tuning (Stiennon et al., 2020; Ouyang et al., 2022; Glaese et al., 2022; Bai et al., 2022), which requires access to the environment or a simulator. In language modeling, Korbak et al. (2023) showed that accounting for the alignment of behavior with the DVF already during imitation learning yields results superior to fine-tuning after-the-fact, however, their approach considers an agent-specific value function. In contrast, we consider learning a policy aligned with a collective value function, and from offline data alone. Credit assignment in multi-agent systems was initially studied in Economics (Shapley, 1953). Subsequently, Shapley Values (Shapley, 1953) and related concepts have been applied in multi-agent reinforcement learn-




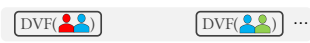

Example Scenario	Dataset characteristics			Example observation dataset of agents  , ...	Shapley Values Applicable	Exchange Values Applicable	Exchange Value Computation
	All Group Sizes	All Observations	Known Agent Identities				
Ideal All group sizes are permitted. All possible combinations of agents (all groups) are observed.	✓	✓	✓		✓	✓	Exact
Group-Limited Only specific group sizes are permitted. A football game has 11 players.	✗	✓	✓		✗	✓	Exact
Low-Data Not all permitted groups are observed. Two football players might never play for the same team.	✗	✗	✓		✗	✓	Estimated (potentially with EV-Clustering)
Degenerate Anonymized with one-time-use IDs. New ID for a player in each game played.	✗	✗	✗		✗	✓	Estimated with EV-Clustering and behaviour information from τ

Figure 2: Overview of different characteristics of real-world datasets, and whether Shapley Values and Exchange Values (EVs) are applicable to compute contributions of individual agents to the DVF.

ing, to distribute rewards among individual agents during the learning process (Chang et al., 2003; Foerster et al., 2018; Nguyen et al., 2018; Wang et al., 2020; Li et al., 2021; Wang et al., 2022). Outside of policy learning, Heuillet et al. (2022) used Shapley Values to analyze agent contributions in multi-agent environments, however this requires privileged access to a simulator, in order to replace agents with randomly-acting agents. In contrast to Shapley Values, the applicability of EVs to all group sizes allows us to omit the need to simulate infeasible coalitions by summing over multiple outcomes or with random-action policies. In contrast to this work, existing work in multi-agent imitation learning typically assumes observations to be generated by optimal agents, as well as simulator access (Le et al., 2017; Song et al., 2018; Yu et al., 2019). Similar to our framework, offline multi-agent reinforcement learning (Jiang & Lu, 2021; Tseng et al., 2022; Tian et al., 2022) involves policy learning from multi-agent demonstrations using offline data alone, however, it assumes a dense reward signal to be given, while the DVF assigns a single score per collective trajectory. In single-agent settings, a large body of work investigates estimating demonstrator expertise to enhance imitation learning (Chen et al., 2021; Zhang et al., 2021; Cao & Sadigh, 2021; Sasaki & Yamashina, 2021; Beliaev et al., 2022; Yang et al., 2021). However, these methods do not translate to the multi-agent setting due to the challenge of credit assignment. To the best of our knowledge, no prior work has considered the problem of imitating multi-agent datasets containing unaligned agents, while ensuring alignment with a collective value function.

3 BACKGROUND AND NOTATION

Markov Game. We consider Markov Games (Littman, 1994), which generalize Markov Decision Processes (MDPs) to multi-agent scenarios. In a Markov Game, agents interact in a common environment. At time step t , each agent (the i th of a total of m agents) takes the action a_i^t and the environment transitions from state s^t to s^{t+1} . A reduced Markov game (without rewards) is then defined by a state space \mathcal{S} ($s^t \in \mathcal{S}$), a distribution of initial states η , the action space \mathcal{A}_i ($a_i^t \in \mathcal{A}_i$) of each agent i , an environment state transition probability $P(s^{t+1}|s^t, a_1, \dots, a_m)$ and the episode length T . We denote this Markov Game as $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, T)$, with collective trajectories $\tau = (s_0, \mathbf{a}_0, \dots, s_T)$.

Set of multi-agent demonstrations generated by many agents. We consider a Markov game \mathcal{M} of m agents and a set of demonstrator agents $N = \{1, \dots, n\}$ where $n \geq m$. The Markov Game further is assumed to be symmetric (we can change the ordering of players without changing the game). The demonstration set \mathcal{D} captures interactions among various groups of agents in \mathcal{M} . Every entry $\mathcal{D}_i = (s_i, \tau_{s_i})$ contains a trajectory τ_{s_i} for a group of agents $s_i \subseteq N$. Notably, τ_{s_i} contains the collective trajectory of all agents in the group s_i .

Shapley Values. We now define the concept of the Shapley Value of an agent Shapley (1953), which is commonly used to evaluate contributions of individual agents to a collective value function

in a characteristic function game. Definition 3.2 below is somewhat unconventional, but can be easily seen to be equivalent to the standard definition.

Definition 3.1 (Characteristic function game). A characteristic function game G is given by a pair (N, v) , where $N = \{1, \dots, n\}$ is a finite, non-empty set of agents and $v : 2^N \rightarrow \mathbb{R}$ is a characteristic function, which maps each group (sometimes also referred to as coalition) $C \subseteq N$ to a real number $v(C)$; it is assumed that $v(\emptyset) = 0$. The number $v(C)$ is referred to as the value of the group C .

Given a characteristic function game $G = (N, v)$, let $\Pi_{N \setminus \{i\}}$ denote the set of all permutations of $N \setminus \{i\}$, i.e., one-to-one mappings from $N \setminus \{i\}$ to itself. For each permutation $\pi \in \Pi_{N \setminus \{i\}}$, we denote by $S_\pi(m)$ the slice of π up until and including position m ; we think of $S_\pi(m)$ as the set of all agents that appear in the first m positions in π (note that $S_\pi(0) = \emptyset$). The marginal contribution of an agent i with respect to a permutation π and a slice m in a game $G = (N, v)$ is given by $\Delta_{m,\pi}^G(i) = v(S_\pi(m) \cup \{i\}) - v(S_\pi(m))$.

This quantity measures the increase in the value of the group when agent i joins them. We can now define the Shapley Value of an agent i : it is simply the agent’s average marginal contribution, where the average is taken over all permutations of $N \setminus \{i\}$ and all slices.

Definition 3.2 (Shapley Value). Given a characteristic function game $G = (N, v)$ with $|N| = n$, the Shapley Value of an agent $i \in N$ is denoted by $SV_i(G)$ and is given by

$$SV_i(G) = 1/n! \cdot \sum_{m=0}^{n-1} \sum_{\pi \in \Pi_{N \setminus \{i\}}} \Delta_{m,\pi}^G(i). \quad (1)$$

Def. 3.2 is important in the context of credit assignment, as a possible solution for distributing collective value to individual agents. It also has several consistency properties (Shapley, 1953).

4 METHODS

Problem setting. Given a dataset \mathcal{D} of trajectories generated by groups of interacting agents and a *Desired Value Function* (DVF), the goal of our paper is to learn an imitation policy for a single agent that is aligned with the DVF. We assume that a fraction of the demonstrator agents’ behavior is undesirable, specifically, their presence in a group results in a significant reduction of the DVF. Further, we assume that the number of demonstrator agents is much larger than the group size of the target scenario.

Overview of Methods section. To evaluate agents’ contributions in games that only permit specific group sizes, we first define the concept of EVs (Def.4.1) for regular characteristic function games (Def. 3.1). We then show that our definition extends naturally to characteristic function games with constraints on permitted group sizes. We finally derive methods to estimate EVs from real-world datasets with limited observations (see Figure 2 for an overview).

4.1 EXCHANGE VALUES TO EVALUATE AGENTS’ INDIVIDUAL CONTRIBUTIONS

Note that each term of the Shapley Value, denoted $\Delta_{m,\pi}^G(i)$, requires computing the difference in values between 2 groups of *different* sizes, with and without an agent i (see Def. 3.2). If we wish to only compare groups with the same size, then a natural alternative is to compute the difference in values when the agent at position m is replaced with agent i :

$$\Gamma_{m,\pi}^G(i) = v(S_\pi(m-1) \cup \{i\}) - v(S_\pi(m)). \quad (2)$$

We call this quantity the *exchange contribution* of i , given a permutation of agents π sliced at position m . It represents the added value of agent i in a group. Importantly it does not require values of groups of different sizes.

We now define the EV analogously to the Shapley Value as the average exchange contribution over all permutations of $N \setminus \{i\}$ and all non-empty slices.

Definition 4.1 (Exchange Value). Given a characteristic function game $G = (N, v)$ with $|N| = n$, the Exchange Value of an agent $i \in N$ is denoted by $EV_i(G)$ and is given by

$$EV_i(G) = ((n-1)! \cdot (n-1))^{-1} \cdot \sum_{m=1}^{n-1} \sum_{\pi \in \Pi_{N \setminus \{i\}}} \Gamma_{m,\pi}^G(i). \quad (3)$$

In words, the EV of an agent can hence be understood as the expected change in value, when substituting the agent with another randomly selected agent, or as comparing the value of all groups that include the agent to that of all groups which do not include the agent (see Step 2 in Figure 1).

Relationship between Shapley Value and Exchange Value. It can be shown that the Exchange Values of all agents can be derived from their Shapley Values by a simple linear transformation: we subtract a fraction of the value of the grand coalition N (group of all agents) and scale the result by $n/n-1$: $EV_i(G) = \frac{n}{n-1}(SV_i(G) - 1/n \cdot v(N))$. The proof proceeds by computing the coefficient of each term $v(C)$, $C \subseteq N$, in summations (1) and (3) (see Appendix A). Hence, the Shapley Value and the Exchange Value order the agents in the same way. Now, recall that the Shapley Value is characterized by four axioms, namely, dummy, efficiency, symmetry and linearity (Shapley, 1953). The latter two are also satisfied by the Exchange Value: if $v(C \cup \{i\}) = v(C \cup \{j\})$ for all $C \subseteq N \setminus \{i, j\}$, we have $EV_i(G) = EV_j(G)$ (symmetry), and if we have two games G_1 and G_2 with characteristic functions v_1 and v_2 over the same set of agents N , then for the combined game $G = (N, v)$ with the characteristic function v given by $v(C) = v_1(C) + v_2(C)$ we have $EV_i(G) = EV_i(G_1) + EV_i(G_2)$ (linearity). The efficiency property of the Shapley Value, i.e., $\sum_{i \in N} SV_i(G) = v(N)$ implies that $\sum_{i \in N} EV_i(G) = 0$. In words, the sum of all agents' EV is zero. The dummy axiom, too, needs to be modified: if an agent i is a dummy, i.e., $v(C \cup \{i\}) = v(C)$ for every $C \subseteq N \setminus \{i\}$ then for the Shapley value we have $SV_i(G) = 0$ and hence $EV_i(G) = -1/n-1 \cdot v(N)$. In each case, the proof follows from the relationship between the Shapley Value and the Exchange Value and the fact that the Shapley Value satisfies these axioms (see Appendix A).

4.1.1 COMPUTING EXCHANGE VALUES IF ONLY CERTAIN GROUP SIZES ARE PERMITTED

For a characteristic function game $\mathcal{G} = (N, v)$ the value function v can be evaluated for each possible group $C \subseteq N$. We now consider the case where the value function v is only defined for groups of certain sizes $m \in M$, i.e. v is only defined for a subset of groups of certain sizes.

Definition 4.2 (Constrained characteristic function game). A constrained characteristic function game \bar{G} is given by a tuple (N, v, M) , where $N = \{1, \dots, n\}$ is a finite, non-empty set of agents, $M \subseteq \{0, \dots, n-1\}$ is a set of feasible group sizes and $v : \{C \in 2^N : |C| \in M\} \rightarrow \mathbb{R}$ is a characteristic function, which maps each group $C \subseteq N$ of size $|C| \in M$ to a real number $v(C)$.

Note that both the Shapley Value and the EV are generally undefined for constrained characteristic function games, as the value function is not defined for groups C of size $|C| \notin M$. The definition of the Shapley Value cannot easily be adapted to constrained characteristic function games, as its computation requires evaluating values of groups of different sizes. In contrast, the definition of the EV can be straightforwardly adapted to constrained characteristic function games by limiting the summation to slices of size $m \in M^+$, where $M^+ = \{m \in M : m > 0\}$. Hence, we define the Constrained EV as the average exchange contribution over all permutations of $N \setminus \{i\}$ and over all slices of size $m \in M^+$.

Definition 4.3 (Constrained Exchange Value). Given a constrained characteristic function game $\bar{G} = (N, v, M)$ with $|N| = n$, the Constrained Exchange Value of an agent $i \in N$ is denoted by $EV_i(\bar{G})$ and is given by $EV_i(\bar{G}) = ((n-1)! \cdot |M^+|)^{-1} \cdot \sum_{m \in M^+} \sum_{\pi \in \Pi_{N \setminus \{i\}}} \Gamma_{m, \pi}^{\bar{G}}(i)$.

We refer to the Constrained EV and EV interchangeably, as they are applicable to different settings. As outlined in Step 2 in Figure 1, the EV of an agent is a comparison of the value of a group that includes the agent and a group that does not include the agent, considering all permitted group sizes.

4.2 ESTIMATING EXCHANGE VALUES FROM LIMITED DATA

The EV assesses the contribution of an individual agent and is applicable under group size limitations in real-world scenarios (see Group-Limited in Figure 2). However, exactly calculating EVs is almost always impossible as real-world datasets likely do not contain observations for all (combinatorially many) possible groups (Low-Data in Figure 2). We first show a sampling-based estimate (Section 4.2) of EVs, which may have a high variance for EVs of agents that are part of a only few trajectories (outcomes). Next we introduce a novel method, EV-Clustering (Section 4.2.1), which clusters agents that behave similarly and can be used to reduce the variance. When datasets are anonymized with one-time-use IDs each demonstrator is only observed as part of one group (see

Degenerate in Figure 2), rendering credit assignment degenerate, as explained in Section 4.2.1; we propose to address this by incorporating low-level behavior data from the trajectories τ . If some groups are not observed, we can achieve an unbiased estimate of the EV by sampling groups uniformly at random. The expected EV is $EV_i(\tilde{G}) = \mathbb{E}_{m \sim U(M^+), \pi \sim U(\Pi_{N \setminus \{i\}})} [\Gamma_{m, \pi}^{\tilde{G}}(i)]$. This expectation converges to the true EV in the limit of infinite samples.

4.2.1 EV-CLUSTERING IDENTIFIES SIMILAR AGENTS

In the case of very few agent observations, the above-introduced sampling estimate has a high variance. One way to reduce the variance is by clustering: if we knew that some agents tend to contribute similarly to the DVF, then by clustering them and estimating one EV per cluster (instead of one EV per agent), each EV estimate will use more samples. Note that, as our focus is on accurately estimating EVs, we do not consider clustering agents by behavior here, as two agents may exhibit distinct behaviors while still contributing similarly to the DVF.

We propose *EV-Clustering*, which clusters agents such that the variance in assigned EVs is maximized across all agents. In Appendix A we show that *EV-Clustering* is equivalent to clustering agents by their *unobserved* individual contribution, under the approximation that the total value of a group is the sum of the participating agents' individual contributions, an assumption frequently made for theoretical analysis (Lundberg & Lee, 2017; Covert & Lee, 2021), as it represents the simplest non-trivial class of cooperative games. Intuitively, if we choose clusters that maximize the EV variance across agents, all clusters' EVs will be maximally distinct. An example of poor clustering is a random partition, which will have low variance and thus very similar EVs across clusters.

Specifically, we assign n agents to $k \leq n$ clusters $K = \{1, \dots, k-1\}$, with individual cluster assignments $\mathbf{u} = \{u_0, \dots, u_{n-1}\}$, where $u_i \in K$. We first combine the observations of all agents within the same cluster by defining a clustered value function $\tilde{v}(C)$ that assigns a value to a group of cluster-centroid agents $C \subseteq K$ by averaging over the combined observations, as $\tilde{v}(C) = 1/\eta \cdot \sum_{m=0}^{n-1} \sum_{\pi \in \Pi_N} v(S_\pi(m)) \cdot \mathbb{1}(\{u_j \mid j \in S_\pi(m)\} = C)$, where η is a normalization constant. The EV of an agent i is then given as $EV_i(\tilde{G})$, with $\tilde{G} = (K, \tilde{v})$, thereby assigning equal EVs to all agents within one cluster.

Definition 4.4 (EV-Clustering). We define the optimal cluster assignments \mathbf{u}^* such that the variance of EVs is maximised:

$$\mathbf{u}^* \in \arg \max_{\mathbf{u}} \text{Var}([EV_0(\tilde{G}), \dots, EV_{n-1}(\tilde{G})]). \quad (4)$$

We show in Appendix B.1 that this objective is equivalent to clustering agents by their unobserved individual contributions, under the approximation of an additive value function.

4.2.2 DEGENERACY OF THE CREDIT ASSIGNMENT PROBLEM FOR FULLY-ANONYMIZED DATA

If two agents are observed only once in the dataset and as part of the same group, equal credit must be assigned to both due to the inability to separate their contributions. Analogously, when all agents are only observed once, credit can only be assigned to groups, resulting in the degenerate scenario that all agents in a group are assigned the same credit (e.g. are assigned equal EV). We solve this by combining low-level behavior information from trajectories τ with EV-Clustering (see Sec. 5.1).

4.3 EXCHANGE VALUE BASED BEHAVIOR CLONING (EV2BC)

Having defined the EV of an individual agent and different methods to estimate it, we now define a variation of Behavior Cloning (Pomerleau, 1991) which takes into account each agent's contribution to the desirability value function (DVF). We refer to this method as *EV2BC*. Essentially, EV2BC imitates only actions of agents that have an EV larger than a tunable threshold parameter.

Definition 4.5 (EV based Behavior Cloning (EV2BC)). For a set of demonstrator agents N , a dataset \mathcal{D} , and a DVF, we define the imitation learning loss for EV2BC as

$$L_{EV2BC}(\theta) = -\sum_{n \in \mathcal{N}} \sum_{(s_i, a_i^n) \in \mathcal{D}} \log(\pi^\theta(a_i^n | s_i)) \cdot \mathbb{1}(EV_n^{DVF} > c) \quad (5)$$

where EV_n^{DVF} denotes the EV of agent n for a DVF and where c is a tunable threshold parameter, that trades off between including data of agents with higher contributions to the DVF and reducing the total amount of training data used.

Table 1: Resulting performance with respect to the DVF for different imitation learning methods in different Starcraft scenarios.

Method	2s3z	3s_vs_5z	6h_vs_8z
BC (Pomerleau, 1991)	13.24 \pm 1.26	12.73 \pm 3.25	9.56 \pm 0.67
Group-BC	17.25 \pm 2.05	12.32 \pm 1.92	10.08 \pm 1.07
EV2BC (Ours)	19.46 \pm 2.98	17.15 \pm 2.13	12.25 \pm 1.55

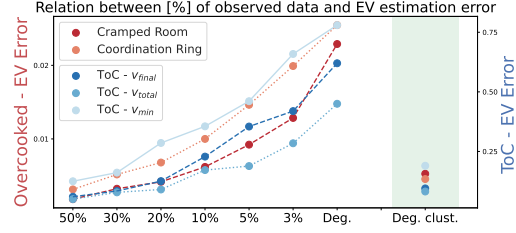


Figure 3: Mean error in estimating EVs with decreasing number of observations. "Deg." refers to the fully anonymized degenerate case. Error decreases significantly if agents are clustered (green-shaded area).

5 EXPERIMENTS

The environments that we consider only permit certain group sizes, hence we use constrained EVs (see Def. 4.3). As the environments are stochastic, we use sampling (see Sec. 4.2) to estimate true EVs. We run all experiments for five random seeds and report mean and standard deviation where applicable. For more implementation details please refer to the Appendix.

In the following experiments, we first evaluate EVs as a measure of an agent’s contribution to a given DVF. We then assess the average estimation error for EVs as the number of observations in the dataset D decreases, and how applying clustering decreases this error. We lastly evaluate the performance of Exchange Value based Behaviour Cloning (EV2BC, see Definition 4.5) for simulated and human datasets and compare to relevant baselines, such as standard Behavior Cloning (Pomerleau, 1991) and Offline Reinforcement Learning (Pan et al., 2022).

Environments. The **Tragedy of the Commons** (Hardin, 1968) (ToC) refers to a situation where multiple individuals deplete a shared resource, and is a social dilemma scenario often studied to model the overexploitation of common resources (Dietz et al., 2003; Ostrom, 2009). We model ToC as a multi-agent environment and consider three DVFs to represent different interpretations of social welfare in the game: the final pool size (v_{final}), the total resources consumed (v_{total}), and the minimum consumption among agents (v_{min}). **Overcooked** (Carroll et al., 2019; Hu et al., 2020; Shih et al., 2022) is a two-player game simulating a cooperative cooking task requiring teamwork and coordination and a common testbed in multi-agent research for studying collaboration. Within Overcooked, we consider the configurations Cramped Room and Coordination Ring (displayed in Figure 4). For each environment configuration, we generate two datasets by simulating agent behaviors using a near-optimal planning algorithm, where we use a parameter λ to determine an agent’s behavior. For $\lambda = 1$ agents act (near)-optimal, for $\lambda = -1$ agents act adversarially. We refer to λ as the agent’s trait, as it acts as a proxy for the agent’s individual contribution to the collective value function. Each demonstration dataset D is generated by $n = 100$ agents, and trajectories τ are of length 400. The adversarial dataset D^{adv} is comprised of 25% adversarial agents with $\lambda = -1$ and 75% (near)-optimal agents with $\lambda = 1$, while for the dataset D^λ agents were uniformly sampled between $\lambda = -1$ and $\lambda = 1$. The D^{human} dataset was collected from humans playing the game (see Carroll et al. (2019)); it is fully anonymized with one-time-use agent identifiers, hence is a degenerate dataset (see Figure 2 bottom row). We consider the standard value function given for Overcooked as the DVF, i.e. the number of soups prepared by both agents over the course of a trajectory. The **StarCraft Multi-Agent Challenge** (Samvelyan et al., 2019) is a cooperative multi-agent environment that is partially observable, involves long-term planning, requires strong coordination, and is heterogeneous, in which we consider the settings 2s3z, 3s_vs_5z and 6h_vs_8z, which involve teams of 3-6 agents. For each, we generate a pool of 200 agents with varying capabilities by extracting policies at different epochs, and from training with different seeds. We generate a dataset that contains simulated trajectories of 100 randomly sampled groups (out of 10^9 possible groups) and use the environment’s ground truth reward function to assign DVF scores according to the collective performance of agents.

Exchange Values assess an agent’s individual contribution to a collective value function. To analyze EVs as a measure for an agent’s individual contribution to a DVF, we consider full datasets that contain demonstrations of all possible groups, which allows us to accurately estimate EVs.

Table 2: Resulting performance with respect to the DVF for different imitation learning methods in the Overcooked environments Cramped Room (top) and Coordination Ring (bottom). In Tragedy of Commons: 12 agents experiment at the top, 120 agents experiment at the bottom.

Imitation method	Overcooked			Overcooked+Fire		Tragedy of Commons		
	\mathcal{D}^λ	\mathcal{D}^{adv}	\mathcal{D}^{human}	\mathcal{D}^λ	\mathcal{D}^{adv}	v_{final}	v_{total}	v_{min}
BC (Pomerleau, 1991)	10.8 \pm 2.14	40.8 \pm 12.7	153.34 \pm 11.5	-13.35 \pm 24.5	-20.12 \pm 18.5	2693.6 \pm 139.1	50.6 \pm	2.4 \pm 0.45
Group-BC	54.2 \pm 5.45	64.8 \pm 7.62	163.34 \pm 6.08	24.89 \pm 16.25	0.9 \pm 13.98	5324.2 \pm 210.8	100.01 \pm 20.08	4.60 \pm 1.01
OMAR (Pan et al., 2022)	6.4 \pm 3.2	25.6 \pm 8.9	12.5 \pm 4.5	5.0 \pm 12.5	-3.4 \pm 12.8	-	-	-
EV2BC (ours)	91.6 \pm 12.07	104.2 \pm 10.28	170.89 \pm 6.8	86.2 \pm 13.02	98.3 \pm 12.48	10576.8 \pm 307.4	342.8 \pm 49.36	44.2 \pm 6.4
BC (Pomerleau, 1991)	15.43 \pm 4.48	10.4 \pm 6.8	104.89 \pm 12.44	-16.45 \pm 15.6	-40 \pm 14.6	2028.8 \pm 60.9	38.9 \pm 10.4	1.8 \pm 0.4
Group-BC	24 \pm 4.69	14.6 \pm 2.48	102.2 \pm 6.19	-8 \pm 8.59	-51.8 \pm 11.4	3400.5 \pm 100.9	77.1 \pm 14.1	3.51 \pm 1.6
OMAR (Pan et al., 2022)	12.43 \pm 3.35	9.5 \pm 3.5	12.4 \pm 6.0	-0.8 \pm 5.4	-1.2 \pm 5.6	-	-	-
EV2BC (ours)	30.2 \pm 6.91	12.4 \pm 2.65	114.89 \pm 5.08	32.64 \pm 7.14	12.5 \pm 4.32	8123.4 \pm 600.8	270.0 \pm 50.0	33.1 \pm 7.1

In ToC, we find that the ordering of agents broadly reflects our intuition, taking more resources negatively impacts the EVs, and agents consuming the average of others have less extreme EVs. The color-coded ordering of agents under different DVFs is shown in Figure 7 in App. C. In Overcooked, we consider the two simulated datasets (\mathcal{D}^{adv} and \mathcal{D}^λ) but not the human dataset, as the individual contribution is unknown for human participants. We find that EVs of individual agents are strongly correlated with their trait parameter λ , which is a proxy for the agent’s individual contribution, and provide a plot that shows the relationship between λ and EV in Figure 5 in App. B).

5.1 ESTIMATING EVs FROM INCOMPLETE DATA

Estimation error for different dataset sizes. We now turn to realistic settings with missing data, where EVs must be estimated (Sec. 4.2). For both ToC and Overcooked, we compute the mean estimation error in EVs if only a fraction of the possible groups is contained in the dataset. As expected, we observe in Fig. 3 that the mean estimation error increases as the fraction of observed groups decreases, with the largest estimation error for fully anonymized datasets (see Fig. 3 – *Deg.*).

Estimating EVs from degenerate datasets. We first use low-level behavior information from the given trajectories τ in D to initialize cluster assignments and then apply EV-Clustering. Specifically, we first create behavior-based cluster assignments by applying k-means clustering to vectors of concatenated action frequencies in frequently-observed states (see Appendix B for details). We then perform EV-Clustering, using the behavior-based cluster assignments to initialize a non-linear constrained optimization solver (SLSQP, Kraft (1988)) and adding a small L_2 loss term that penalizes solutions deviating from the behavior-based clusters. We observe in Figure 3 that this results in a significant decrease in the estimation error of EVs (see – *Deg. clustered*). Generally, EV-clustering is preferable to behavior clustering, as two agents may have equal contributions to the DVF while showing different behaviors; only in cases where only few outcomes per agent are observed is it necessary to also use behavior clustering. In the ablation study in Appendix B.1 we investigate both methods, finding that both behavior-clustering and EV-Clustering are significant, while behavior clustering is more robust in low-data scenarios (as it incorporates all low-level information contained within a trajectory, while EV-Clustering only considers final outcomes).

Estimating EVs from degenerate human datasets in Overcooked. In contrast to the simulated datasets, no estimation error in EVs can be computed for the human-generated datasets as the ground truth EVs are unknown. Also, no latent trait λ that indicates how well a human participant is aligned with the DVF is known. However, to evaluate the quality of the estimated EVs for the human dataset, we use the keystrokes per second of an agent as a proxy for its individual contribution, which we refer to as λ_{human} . We here follow Carroll et al. (2019), which found that this proxy is highly correlated with overall performance. We estimate EVs for human participants as before, relying both on behavior information and DVF scores. We evaluate the quality of computed EVs as the inverse of the within-cluster variance in λ_{human} . Relative to the average within-cluster variance under random cluster assignments, we find that behavior-based clustering results in a reduction of 16% and 25% in Cramped Room and Coordination ring, respectively, while EV-Clustering reduces the within-cluster variance by another 34% and 48% percent, respectively. These findings validate that maximizing variance in EVs allows clustering agents by their individual contributions.

5.2 IMITATING DESIRED BEHAVIOR BY UTILIZING EVs

We now evaluate EV2BC in both domains, where we set the threshold parameter such that in ToC only agents with EVs above the 90th percentile are imitated, and in Overcooked above the 50th

percentile. We chose these values because training data in Overcooked is more scarce. We replicate the stochastic EV2BC policy for each agent in the environment and evaluate the achieved collective DVF score. As baselines, we consider (1) *BC*, where Behavior Cloning (Pomerleau, 1991) is done with the full dataset without correcting for EVs, (2) offline multi-agent reinforcement *OMAR* (Pan et al., 2022) with the reward at the last timestep set to the DVF’s score for a given trajectory (no per-step reward is given by the DVF) and (3) *Group BC*, for which only collective trajectories with a DVF score larger than the relevant percentile are included. While EV2BC is based on individual agents’ contributions, this last baseline imitates data based on group outcomes. For instance, if a collective trajectory includes two aligned agents and one unaligned agent, the latter baseline is likely to imitate all three agents. In contrast, our approach would only imitate the two aligned agents.

ToC results. We imitate datasets of 12 agents and 120 agents, with group sizes of 3 and 10 respectively, evaluating performance for each of the three DVFs defined for the ToC environment. We do not consider the OMAR baseline as policies are not learned but rule-based. Table 2 demonstrates that EV2BC outperforms the baselines by a large margin, indicating that considering individual agents’ EVs to a given DVF leads to significantly improved performance.

Overcooked results. We now consider all datasets D^{adv} , D^{λ} and D^{human} in both Overcooked environments. Note that in the standard Overcooked environment, an adversarial agent is limited to *blocking* the other agent, while in many real-world environments, adversaries are likely to be capable of more diverse (and possibly severe) actions. We evaluate the performance achieved by agents with respect to the DVF (in this case the environments value function of maximizing the number of soups) when trained with different imitation learning approaches on the different datasets. We evaluate performance for the fully-anonymized datasets, but also consider datasets with more data in Table 4 in the Appendix, for which we find an even larger performance gap. EVs are computed as detailed in Section 5.1. Table 2 shows that EV2BC clearly outperforms the baseline approaches in both environment configurations, with the margin being more significant in the Overcooked+Fire environments where adversarial agents can take more powerful actions. We further note that EV2BC significantly outperforms baseline approaches on the datasets of human-generated behavior, for which EVs were estimated from a fully-anonymized real-world dataset. This demonstrates that BC on datasets containing unaligned behavior carries risk of learning wrong behavior, but it can be alleviated by weighting the samples using estimated EVs.

Starcraft Results. We observe in Table 1 that EV2BC outperforms the baselines by a substantial margin, underlining the applicability of our method to larger and more complex settings. Note that the OMAR baseline, which is implemented as offline MARL with the DVF as the final-timestep reward, did substantially worse than BC.

6 CONCLUSION

Our work presents a method for training AI agents from diverse datasets of human interactions while ensuring that the resulting policy is aligned with a given desirability value function. However, it must be noted that quantifying this value function is an active research area.

Shapley Values and Exchange Values estimate the alignment of an individual with a group value function (which must be prescribed separately), and as such can be misused if they are included in a larger system that is used to judge those individuals in any way. Discrimination of individuals based on protected attributes is generally unlawful, and care must be taken to avoid any discrimination by automated means. We demonstrated a novel positive use of these methods by using them to train aligned (beneficial) agents, that do not imitate negative behaviors in a dataset. We expect that the benefits of addressing the problem of unsafe behavior by AI agents outweigh the downsides of misuse of Shapley Values and Exchange Values, which are covered by existing laws.

Future work may address the assumption that individual agents behave similarly across multiple trajectories and develop methods for a more fine-grained assessment of desired behavior. Additionally, exploring how our framework can more effectively utilize data on undesired behavior is an interesting avenue for further investigation, e.g., developing policies that are constrained to not taking undesirable actions. Lastly, future work may investigate applications to real-world domains, such as multi-agent autonomy scenarios.

389 **Reproducibility.** To ensure the reproducibility of our work, we will publish the source code with
 390 the camera-ready version of this work. We provide detailed overviews for all steps of the exper-
 391 imental evaluation in the Appendix, where we also link to the publicly available code repositories
 392 that our work used. We further provide information about computational complexity at the end of
 393 the Appendix.

394 REFERENCES

- 395 Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav
 396 Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with reinforcement
 397 learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- 398 Mark Beliaev, Andy Shih, Stefano Ermon, Dorsa Sadigh, and Ramtin Pedarsani. Imitation learning by estimat-
 399 ing expertise of demonstrators. In *International Conference on Machine Learning*, pp. 1732–1748. PMLR,
 400 2022.
- 401 Zhangjie Cao and Dorsa Sadigh. Learning from imperfect demonstrations from agents with varying dynamics.
 402 *IEEE Robotics and Automation Letters*, 6(3):5231–5238, 2021.
- 403 Micah Carroll, Rohin Shah, Mark K Ho, Tom Griffiths, Sanjit Seshia, Pieter Abbeel, and Anca Dragan. On
 404 the utility of learning about humans for human-ai coordination. *Advances in neural information processing*
 405 *systems*, 32, 2019.
- 406 Yu-Han Chang, Tracey Ho, and Leslie Kaelbling. All learning is local: Multi-agent learning in global reward
 407 games. *Advances in neural information processing systems*, 16, 2003.
- 408 Letian Chen, Rohan Paleja, and Matthew Gombolay. Learning from suboptimal demonstration via self-
 409 supervised reward regression. In *Conference on robot learning*, pp. 1262–1277. PMLR, 2021.
- 410 Ian Covert and Su-In Lee. Improving kernelshap: Practical shapley value estimation using linear regression. In
 411 *International Conference on Artificial Intelligence and Statistics*, pp. 3457–3465. PMLR, 2021.
- 412 Thomas Dietz, Elinor Ostrom, and Paul C Stern. The struggle to govern the commons. *science*, 302(5652):
 413 1907–1912, 2003.
- 414 Meta Fundamental AI Research Diplomacy Team (FAIR)[†], Anton Bakhtin, Noam Brown, Emily Dinan,
 415 Gabriele Farina, Colin Flaherty, Daniel Fried, Andrew Goff, Jonathan Gray, Hengyuan Hu, et al. Human-
 416 level play in the game of diplomacy by combining language models with strategic reasoning. *Science*, 378
 417 (6624):1067–1074, 2022.
- 418 Jakob Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson. Counter-
 419 factual multi-agent policy gradients. In *Proceedings of the AAAI conference on artificial intelligence*, vol-
 420 ume 32, 2018.
- 421 Amelia Glaese, Nat McAleese, Maja Trebacz, John Aslanides, Vlad Firoiu, Timo Ewalds, Maribeth Rauh,
 422 Laura Weidinger, Martin Chadwick, Phoebe Thacker, et al. Improving alignment of dialogue agents via
 423 targeted human judgements. *arXiv preprint arXiv:2209.14375*, 2022.
- 424 Garrett Hardin. The tragedy of the commons: the population problem has no technical solution; it requires a
 425 fundamental extension in morality. *science*, 162(3859):1243–1248, 1968.
- 426 Jerry Zhi-Yang He, Zackory Erickson, Daniel S Brown, Aditi Raghunathan, and Anca Dragan. Learning
 427 representations that enable generalization in assistive tasks. In *Conference on Robot Learning*, pp. 2105–
 428 2114. PMLR, 2023.
- 429 Alexandre Heuillet, Fabien Couthouis, and Natalia Díaz-Rodríguez. Collective explainable ai: Explaining co-
 430 operative strategies and agent contribution in multiagent reinforcement learning with Shapley values. *IEEE*
 431 *Computational Intelligence Magazine*, 17(1):59–71, 2022.
- 432 Hengyuan Hu, Adam Lerer, Alex Peysakhovich, and Jakob Foerster. “other-play” for zero-shot coordination.
 433 In *International Conference on Machine Learning*, pp. 4399–4410. PMLR, 2020.
- 434 Jiechuan Jiang and Zongqing Lu. Offline decentralized multi-agent reinforcement learning. *arXiv preprint*
 435 *arXiv:2108.01832*, 2021.
- 436 Tomasz Korbak, Kejian Shi, Angelica Chen, Rasika Bhalerao, Christopher L Buckley, Jason Phang, Samuel R
 437 Bowman, and Ethan Perez. Pretraining language models with human preferences. *arXiv preprint*
 438 *arXiv:2302.08582*, 2023.

- Dieter Kraft. A software package for sequential quadratic programming. *Forschungsbericht- Deutsche Forschungs- und Versuchsanstalt für Luft- und Raumfahrt*, 1988.
- Hoang M Le, Yisong Yue, Peter Carr, and Patrick Lucey. Coordinated multi-agent imitation learning. In *International Conference on Machine Learning*, pp. 1995–2003. PMLR, 2017.
- Jiahui Li, Kun Kuang, Baoxiang Wang, Furui Liu, Long Chen, Fei Wu, and Jun Xiao. Shapley counterfactual credits for multi-agent reinforcement learning. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pp. 934–942, 2021.
- Michael L. Littman. Markov games as a framework for multi-agent reinforcement learning. In *Machine Learning Proceedings 1994*. 1994. doi: 10.1016/b978-1-55860-335-6.50027-1.
- Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *Advances in neural information processing systems*, 30, 2017.
- Duc Thien Nguyen, Akshat Kumar, and Hoong Chuin Lau. Credit assignment for collective multiagent rl with global rewards. *Advances in neural information processing systems*, 31, 2018.
- Elinor Ostrom. A general framework for analyzing sustainability of social-ecological systems. *Science*, 325(5939):419–422, 2009.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.
- Ling Pan, Longbo Huang, Tengyu Ma, and Huazhe Xu. Plan better amid conservatism: Offline multi-agent reinforcement learning with actor rectification. In *International Conference on Machine Learning*, pp. 17221–17237. PMLR, 2022.
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- Dean A. Pomerleau. Efficient Training of Artificial Neural Networks for Autonomous Navigation. *Neural Computation*, 3(1), 1991. ISSN 0899-7667. doi: 10.1162/neco.1991.3.1.88.
- Mikayel Samvelyan, Tabish Rashid, Christian Schroeder De Witt, Gregory Farquhar, Nantas Nardelli, Tim GJ Rudner, Chia-Man Hung, Philip HS Torr, Jakob Foerster, and Shimon Whiteson. The starcraft multi-agent challenge. *arXiv preprint arXiv:1902.04043*, 2019.
- Fumihiko Sasaki and Ryota Yamashina. Behavioral cloning from noisy demonstrations. In *International Conference on Learning Representations*, 2021.
- Lloyd Shapley. A value for n -person games. *Contributions to the Theory of Games*, pp. 307–317, 1953.
- Andy Shih, Stefano Ermon, and Dorsa Sadigh. Conditional imitation learning for multi-agent games. In *2022 17th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*, pp. 166–175. IEEE, 2022.
- Jiaming Song, Hongyu Ren, Dorsa Sadigh, and Stefano Ermon. Multi-agent generative adversarial imitation learning. *Advances in neural information processing systems*, 31, 2018.
- Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. *Advances in Neural Information Processing Systems*, 33:3008–3021, 2020.
- Robert Thorndike. Who belongs in the family? *Psychometrika*, 18(4):267–276, 1953.
- Qi Tian, Kun Kuang, Furui Liu, and Baoxiang Wang. Learning from good trajectories in offline multi-agent reinforcement learning. *arXiv preprint arXiv:2211.15612*, 2022.
- Wei-Cheng Tseng, Tsun-Hsuan Johnson Wang, Yen-Chen Lin, and Phillip Isola. Offline multi-agent reinforcement learning with knowledge distillation. *Advances in Neural Information Processing Systems*, 35: 226–237, 2022.
- Jianhong Wang, Yuan Zhang, Tae-Kyun Kim, and Yunjie Gu. Shapley q-value: A local reward approach to solve global reward games. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 7285–7292, 2020.

- 488 Jianhong Wang, Yuan Zhang, Yunjie Gu, and Tae-Kyun Kim. Shaq: Incorporating shapley value theory into
489 multi-agent q-learning. *Advances in Neural Information Processing Systems*, 35:5941–5954, 2022.
- 490 Jing Xu, Da Ju, Margaret Li, Y-Lan Boureau, Jason Weston, and Emily Dinan. Recipes for safety in open-
491 domain chatbots. *arXiv preprint arXiv:2010.07079*, 2020.
- 492 Mengjiao Yang, Sergey Levine, and Ofir Nachum. Trail: Near-optimal imitation learning with suboptimal data.
493 *arXiv preprint arXiv:2110.14770*, 2021.
- 494 Lantao Yu, Jiaming Song, and Stefano Ermon. Multi-agent adversarial inverse reinforcement learning. In
495 *International Conference on Machine Learning*, pp. 7194–7201. PMLR, 2019.
- 496 Songyuan Zhang, Zhangjie Cao, Dorsa Sadigh, and Yanan Sui. Confidence-aware imitation learning from
497 demonstrations with varying optimality. *Advances in Neural Information Processing Systems*, 34:12340–
498 12350, 2021.

A APPENDIX TO METHODS

A.1 AXIOMATIC PROPERTIES OF THE EXCHANGE VALUE AND ITS RELATIONSHIP WITH THE SHAPLEY VALUE

Fix a characteristic function game G with a set of players N . It is well-known that the Shapley Value satisfies the following axioms (Shapley, 1953):

- (1) Dummy: if an agent i satisfies $v(C \cup \{i\}) = v(C)$ for all $C \subseteq N \setminus \{i\}$ then $SV_i(G) = 0$;
- (2) Efficiency: the sum of all agents' Shapley Values equals to the value of the grand coalition, i.e., $\sum_{i \in N} SV_i(G) = v(N)$;
- (3) Symmetry: for every pair of distinct agents $i, j \in N$ with $v(C \cup \{i\}) = v(C \cup \{j\})$ for all $C \subseteq N \setminus \{i, j\}$ we have $SV_i(G) = SV_j(G)$;
- (4) Linearity: for any pair of games $G_1 = (N, v_1)$ and $G_2 = (N, v_2)$ with the same set of agents N , the game $G = (N, v)$ whose characteristic function v is given by $v(C) = v_1(C) + v_2(C)$ for all $C \subseteq N$ satisfies $SV_i(G) = SV_i(G_1) + SV_i(G_2)$ for all $i \in N$.

Indeed, the Shapley Value is the only value for characteristic function games that satisfies these axioms (Shapley, 1953). It is then natural to ask which of these axioms (or their variants) are satisfied by the Exchange Value.

To answer this question, we first establish a relationship between the Shapley Value and the Exchange Value.

Proposition A.1. *For any characteristic function game $G = (N, v)$ and every agent $i \in N$ we have*

$$EV_i(G) = \frac{n}{n-1} \left(SV_i(G) - \frac{1}{n} \cdot v(N) \right). \quad (6)$$

Proof. Fix an agent i and consider an arbitrary non-empty coalition $C \subsetneq N \setminus \{i\}$.

In the expression for the Shapley Value of i the coefficient of $v(C)$ is

$$-\frac{1}{n!}(|C|)!(n-1-|C|)! :$$

we subtract the fraction of permutations of N where the agents in C appear in the first $|C|$ positions, followed by i . By the same argument, the coefficient of $v(C \cup \{i\})$ is

$$\frac{1}{n!}(|C|)!(n-1-|C|)!.$$

Similarly, in the expression for the Exchange Value of i the coefficient of $v(C)$ is

$$-\frac{1}{(n-1)!(n-1)}(|C|)!(n-1-|C|)! :$$

each permutation of $N \setminus \{i\}$ where agents in C appear in the first $|C|$ positions contributes with coefficient $-\frac{1}{(n-1)!(n-1)}$. By the same argument, the coefficient of $v(C \cup \{i\})$ is

$$\frac{1}{(n-1)!(n-1)}(|C|)!(n-1-|C|)!$$

Now, if $C = N \setminus \{i\}$, in the expression for $SV_i(G)$ the coefficient of $v(C)$ is $-\frac{1}{n}$ and the coefficient of $v(C \cup \{i\}) = v(N)$ is $\frac{1}{n}$. In contrast, in the expression for $EV_i(G)$ the coefficient of $v(C)$ is $-\frac{1}{n-1}$: for each of the $(n-1)!$ permutations of $N \setminus \{i\}$ we subtract $v(C)$ with coefficient $\frac{1}{(n-1)!(n-1)}$ when we replace the last agent in that permutation by i . On the other hand, $v(N)$ never appears.

It follows that, for every coalition $C \subsetneq N$, if the value $v(C)$ appears in the expression for $SH_i(G)$ with weight ω then it appears in the expression for $EV_i(G)$ with weight $\frac{n}{n-1} \cdot \omega$. Hence

$$EV_i(G) = \frac{n}{n-1} \left(SH_i(G) - \frac{1}{n} \cdot v(N) \right)$$

522

□

Example A.2. Consider a characteristic function game $G = (N, v)$, where $N = \{1, 2\}$ and v is given by $v(\{1\}) = 2$, $v(\{2\}) = 4$, $v(\{1, 2\}) = 10$. We have

$$SH_1(G) = (2 + (10 - 4))/2 = 4, \quad SH_2(G) = (4 + (10 - 2))/2 = 6$$

and

$$EV_1(G) = 2 - 4 = -2, \quad EV_2(G) = 4 - 2 = 2.$$

523 Note that $EV_1(G) = 2(SH_1(G) - \frac{1}{2}v(N))$, $EV_2(G) = 2(SH_2(G) - \frac{1}{2}v(N))$.

524 We can use Proposition A.1 to show that the Exchange Value satisfies two of the axioms satisfied by the Shapley
525 Value, namely, linearity and symmetry.

526 **Proposition A.3.** *The Exchange Value satisfies symmetry and linearity axioms.*

Proof. For the symmetry axiom, fix a characteristic function game $G = (N, v)$ and consider two agents $i, j \in N$ with $v(C \cup \{i\}) = v(C \cup \{j\})$ for all $C \subseteq N \setminus \{i, j\}$. We have

$$EV_i(G) = \frac{n}{n-1} \left(SV_i(G) - \frac{1}{n} \cdot v(N) \right) = \frac{n}{n-1} \left(SV_j(G) - \frac{1}{n} \cdot v(N) \right) = EV_j(G),$$

527 where the first and the third equality follow by Proposition A.1, and the second equality follows because the
528 Shapley Value satisfies symmetry.

529 For the linearity axiom, consider a pair of games $G_1 = (N, v_1)$ and $G_2 = (N, v_2)$ with the same set of agents
530 N and the game $G = (N, v)$ whose characteristic function v is given by $v(C) = v_1(C) + v_2(C)$ for all
531 $C \subseteq N$. Fix an agent $i \in N$. We have

$$\begin{aligned} EV_i(G) &= \frac{n}{n-1} \left(SV_i(G) - \frac{1}{n} \cdot (v_1(N) + v_2(N)) \right) \\ &= \frac{n}{n-1} \left(SV_i(G_1) - \frac{1}{n} \cdot v_1(N) \right) + \frac{n}{n-1} \left(SV_i(G_2) - \frac{1}{n} \cdot v_2(N) \right) \\ &= EV_i(G_1) + EV_i(G_2). \end{aligned}$$

532 Again, the first and the third equality follow by Proposition A.1, and the second equality follows because the
533 Shapley Value satisfies linearity. \square

534 While the Exchange Value does not satisfy the dummy axiom or the efficiency axiom, it satisfies appropriately
535 modified versions of these axioms.

536 **Proposition A.4.** *For every characteristic function game G it holds that $\sum_{i \in N} EV_i(G) = 0$. Moreover, if i is
537 a dummy agent, i.e., $v(C \cup \{i\}) = v(C)$ for all $C \subseteq N \setminus \{i\}$ then $EV_i(G) = -\frac{v(N)}{n-1}$.*

538 *Proof.* We have

$$\begin{aligned} \sum_{i \in N} EV_i(G) &= \sum_{i \in N} \frac{n}{n-1} \left(SV_i(G) - \frac{1}{n} \cdot v(N) \right) = \sum_{i \in N} \frac{n}{n-1} SV_i(G) - \frac{n}{n-1} \cdot v(N) \\ &= \frac{n}{n-1} \cdot v(N) - \frac{n}{n-1} \cdot v(N) = 0, \end{aligned}$$

539 where we use Proposition A.1 and the fact that the Shapley Value satisfies the efficiency axiom.

Now, fix a dummy agent i . We have

$$EV_i(G) = \frac{n}{n-1} \left(SV_i(G) - \frac{1}{n} \cdot v(N) \right) = -\frac{1}{n-1} \cdot v(N);$$

540 again, we use Proposition A.1 and the fact that the Shapley Value satisfies the dummy axiom. \square

541 A.2 DERIVATION OF CLUSTERING OBJECTIVE STATED IN EQ. 4

542 **Inessential games and EVs.** The assumption of an inessential game is commonly made to compute Shap-
543 ley Values more efficiently². In an inessential game, the value of a group is given by the sum of the individual
544 contributions of its members, denoted as $v(C) = \sum_{i \in C} v_i$, where v_i is an individual agent's unobserved
545 contribution v_i . The EV (see Definition 4.1) of an individual agent i in an inessential game is given as

$$EV_i(G) = v_i - \frac{1}{|N|-1} \cdot \sum_{j \in N \setminus \{i\}} v_j = (1 + \frac{1}{|N|-1}) \cdot v_i - \frac{1}{|N|-1} \cdot \sum_{j \in N} v_j,$$

546 This expression represents the difference between the individual contribution of agent i , v_i , and the average
547 individual contribution of all other agents. The second term is independent of i and remains constant across all
548 agents.

²see, e.g., Covert, I. and Lee, S.I., 2020. Improving kernelshap: Practical shapley value estimation via linear regression. arXiv preprint arXiv:2012.01536.

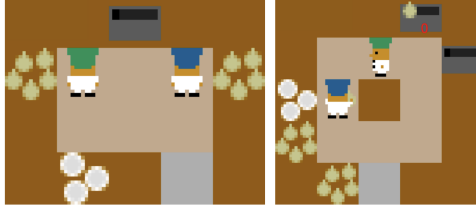


Figure 4: In the Overcooked environments Cramped Room (left) and Coordination Ring (right), agents must cooperate to cook and deliver as many soups as possible within a given time.

Derivation of equivalent clustering objective. We now consider the optimization problem defined by Equation 4, which defines optimal cluster assignments \mathbf{u}^* such that the variance in EVs is maximised

$$\mathbf{u}^* \in \arg \max_{\mathbf{u}} \text{Var}([EV_0(\tilde{G}), \dots, EV_{n-1}(\tilde{G})]).$$

Further, the clustered value function is defined as

$$\tilde{v}(C) = 1/\eta \cdot \sum_{m=0}^{n-1} \sum_{\pi \in \Pi_N} v(S_{\pi}(m)) \cdot \mathbb{1}(\{u_j \mid j \in S_{\pi}(m)\} = C),$$

549 where the normalisation constant is defined as $\eta = \sum_{m=0}^{n-1} \sum_{\pi \in \Pi_N} \mathbb{1}(\{u_j \mid j \in S_{\pi}(m)\} = C)$. We denote by
 550 k_i the individual contribution of the agent that represents the agents in cluster i . The value k_i is defined as the
 551 average individual contribution of all agents assigned to the cluster, i.e. $k_i = 1/\epsilon \cdot \sum_{j \in N} v_j \cdot \mathbb{1}(\mathbf{u}(i) = \mathbf{u}(j))$.
 552 Here, the normalization constant is given as $\epsilon = \sum_{j \in N} \mathbb{1}(\mathbf{u}(i) = \mathbf{u}(j))$.

Using the concept of the clustered value function \tilde{v} , we can express the EV for all agents assigned cluster i as

$$EV_i(\tilde{G}) = (1 + 1/|K|-1) \cdot k_i - 1/|K|-1 \cdot \sum_{j \in K} k_j.$$

553 The second term, which is cluster-independent, can be omitted when computing the variance
 554 $\text{Var}([EV_0(\tilde{G}), \dots, EV_{n-1}(\tilde{G})])$, as the variance is agnostic to a shift in the data distribution. We will omit
 555 the scaling factor $(1 + 1/|K|-1)$ from here onwards.

Let n_j denote the number of agents assigned to cluster $j \in K$, with $\sum_{i=0}^{K-1} n_i = N$. By simplifying Equation 4, we obtain:

$$\text{Var}([EV_0(\tilde{G}), \dots, EV_{n-1}(\tilde{G})]) = \sum_{i=0}^{K-1} n_i \cdot \left(k_i - \sum_{j=0}^{K-1} n_j \cdot k_j / N \right)^2.$$

This allows us to express the objective stated in Equation 4 as

$$\mathbf{u}^* \in \arg \max_{\mathbf{u}} \text{Var}([k_0, \dots, k_{n-1}]).$$

556 The objective stated in Equation 4 is therefore equivalent to assigning agents to clusters such that the variance
 557 in cluster centroids (centroids computed as the mean of the unobserved individual contributions v_i of all agents
 558 assigned to a given cluster) is maximized.

Table 3: Dataset statistics in Overcooked.

Imitation method	Cramped Room D^{λ}	Coordination Ring D^{λ}	Cramped Room D^{adv}	Coordination Ring D^{adv}
Minimum	0	0	0	0
Mean	20.6 ± 33.58	12 ± 19.39	16.91 ± 40.64	3 ± 11.15
Maximum	150	80	160	80

559 B OVERCOOKED EXPERIMENTS

560 We generate the simulated datasets using the planning algorithm given in (Carroll et al., 2019)³. To be able to
 561 simulate agents with different behaviors (from adversarial to optimal), we first introduce a latent trait parameter,

³https://github.com/HumanCompatibleAI/overcooked_ai

λ , which determines the level of adversarial or optimal actions for a given agent. A value of $\lambda = 1$ represents a policy that always chose the best action with certainty. As λ decreases, agents are more likely to select non-optimal actions. For $\lambda < 0$, we invert the cost function to create agents with adversarial behavior. Notably, we assign a high cost (or low cost when inverted) to occupying the cell next to the counter in the Overcooked environment. Occupying the cell next to the counter enables adversarial agents to block other agents in the execution of their tasks.

For human gameplay datasets, we utilized the raw versions of the Overcooked datasets.⁴ These datasets were used as-is, without manual pre-filtering.

We introduce an additional modified version of the Overcooked environment in which agents can take an additional action that lights the kitchen on fire with a predefined probability, resulting in an episode reward of -200 ; we refer to this environment as *Overcooked+Fire* and evaluate on equivalently created datasets D^{adv} and D^λ .

EVs. To estimate agents’ EVs according to Section 4.2, we used either the full set of all possible groups or a fraction of it (see Figure 3 for the relationship between dataset size and EV estimation error). For each observed group, we conducted 10 rollouts in the environment and calculated the average score across these rollouts to account for stochasticity in the environment.

Imitation learning. For EV2BC, BC, and group-BC, we used the implementation of Behavior Cloning in Overcooked as given by the authors of (Carroll et al., 2019)⁵. We implement the offline multi-agent reinforcement learning method OMAR (Pan et al., 2022) using the author’s implementation.⁶ For the OMAR baseline, we set the reward at the last timestep to the DVF’s score for a given trajectory, as our work assumes that no per-step reward signal is given, in contrast to the standard offline-RL framework. We conducted a hyperparameter sweep for the following parameters: learning rate with options $\{0.01, 0.001, 0.0001\}$, Omar-coe with options $\{0.1, 1, 10\}$, Omar-iters with options $\{1, 3, 10\}$, and Omar-sigma with options $\{1, 2, 3\}$. The best-performing parameters were selected based on the evaluation results.

Implementation of Overcooked+Fire. We introduce an additional adversarial action, “light kitchen on fire,” to the environment. To account for this action in the planning algorithm, we assign it the highest possible cost. Taking this action had a 50% chance of resulting in an episode return of -200 , regardless of the other agent’s performance.

B.1 CLUSTERING OF AGENTS IN OVERCOOKED

Behavior clustering. The behavior clustering process in the Overcooked environment involves the following steps. Initially, we identify the 200 states that are most frequently visited by all agents in the given set of observations. As the action space in Overcooked is relatively small (≤ 7 actions), we calculate the empirical action distribution for each state for every agent. These 200 action distributions are then concatenated to form a behavior embedding for each agent. To reduce the dimensionality of the embedding, we apply Principal Component Analysis (PCA), transforming the initial embedding space into three dimensions. Subsequently, we employ the k-means clustering algorithm to assign agents to behavior clusters. The number of clusters (7 for Overcooked) is determined using the ELBOW method (Thorndike, 1953), while linear kernels are utilized for both PCA and k-means. It is noteworthy that the results are found to be relatively insensitive to the parameters used in the dimensionality reduction and clustering steps, thus standard implementations are employed for both methods (Pedregosa et al., 2011). Importantly, this clustering procedure focuses exclusively on the observed behavior of agents, specifically the actions taken in specific states, and is independent of the scores assigned to trajectories by the DVF.

EV-Clustering. In contrast to behavior clustering, EV-Clustering (see Section 4.2.1) focuses solely on the scores assigned to trajectories by the DVF and disregards agent behavior. The objective of variance clustering is to maximize the variance in assigned EVs, as stated in Equation 4. To optimize this objective, we utilize the SLSQP non-linear constrained optimization introduced by Kraft (1988).

We use soft cluster assignments and enforce constraints to ensure that the total probability is equal to one for each agent. The solver is initialized with a uniform distribution and runs until convergence or for a maximum of 100 steps. Given that the optimization problem may have local minima, we perform 500 random initializations and optimizations, selecting the solution with the lowest loss (i.e. the highest variance in assigned EVs).

⁴https://github.com/HumanCompatibleAI/human_aware_rl/tree/master/human_aware_rl/data/human/anonymized

⁵https://github.com/HumanCompatibleAI/overcooked_ai/tree/master/src/human_aware_rl/imitation

⁶<https://github.com/ling-pan/OMAR>

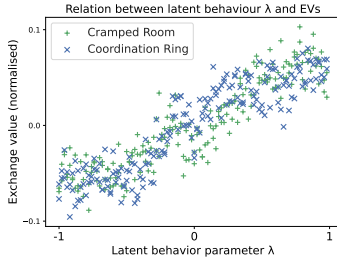


Figure 5: Relation between an agent’s trait λ and its EV in Overcooked.

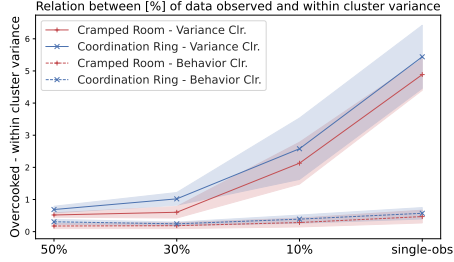


Figure 6: Within-cluster variance in relation to fraction of observations for simulated data in Cramped Room and Coordination Ring (Overcooked). Two clustering methods shown (Behavior clustering and Variance Clustering). In the case of random cluster assignments, the within-cluster variance is 5.11 ± 0.11 , while under optimal cluster assignments, the variance is 0.156. See section B.1 for discussion.

Combining Behavior Clustering and EV Clustering. As described in Sections 4.2.2 and 5.1, behavior clustering (which utilizes behavior information but disregards DVF scores) and variance clustering (which utilizes DVF scores but disregards behavior information) are combined to estimate EVs for degenerate datasets. We initialize the SLSQP solver with the cluster assignments obtained from behavior clustering and introduce a small loss term in the objective function of Equation 4. This additional loss term, weighted by 0.1 (selected in a small sensitivity analysis), penalizes deviations from the behavior clusters. Similar to before, we perform 500 iterations while introducing a small amount of noise to the initial cluster assignments at each step. The solution with the highest variance in assigned EVs is then selected.

Ablation study. We present an ablation study to examine the impact of different components in the clustering approach discussed in Section 5.1. We proposed two sequential clustering methods: behavior clustering and variance clustering. This ablation study investigates the performance of both clustering steps when performed independently, also under the consideration of the fraction of the data that is observed. We assess performance as the within-cluster variance in the unobserved agent-specific latent trait variable λ , where lower within-cluster variance indicates higher performance. It is important to note that λ is solely used for evaluating the clustering steps and not utilized during the clustering process. The results of the ablation study are depicted in Figure 6.

We first discuss EV-Clustering. EV-Clustering as introduced in Section 4 generally leads to a significant decrease in within-cluster variance in the unobserved variable λ . More specifically, the proposed variance clustering approach (when 50% of data is observed), results in a $\sim 89\%$ reduction of the within-cluster variance in λ , which validates the approach of clustering agents by their unobserved individual contributions by maximizing the variance in estimated EVs. However, we observe in Figure 6 that, as the fraction of observed data decreases, the within-cluster variance increases, indicating a decrease in the quality of clustering. The highest within-cluster variance is observed when using only a single observation (‘single-obs’), which corresponds to a fully-anonymized dataset. This finding is consistent with the fact that a fully-anonymized dataset presents a degenerate credit assignment problem, as discussed in Section 4.2.2.

We now discuss behavior clustering. Figure 6 shows that behavior clustering generally results in a very low within-cluster variance. However, it is important to note that these results may not directly translate to real-world data, as the ablation study uses simulated trajectories. Note that such an ablation study cannot be conducted for the given real-world human datasets, as these are fully anonymized. In Section 5.1, we demonstrate that behavior clustering alone may not be sufficient for fully-anonymized real-world human datasets. Instead, a combination of both behavior clustering and variance clustering yields superior results.

Additional results for non-fully anonymized datasets. While the results presented in Section 5.2 were obtained for fully-anonymized datasets, we ran the same evaluations also for the simulated datasets where 30% possible groups are observed. As it can be seen in Table 4, EV2BC outperforms baseline approaches by an even larger margin.

Table 4: Results for 30% of coalitions observed in Overcooked

Imitation method	Cramped Room D^λ	Coordination Ring D^λ	Cramped Room D^{adv}	Coordination Ring D^{adv}
BC	12.6 \pm 3.34	18.13 \pm 6.21	31.7 \pm 8.96	14.21 \pm 3.78
Group-BC	64.33 \pm 6.1	29.4 \pm 7.01	75.7 \pm 13.98	16.8 \pm 5.66
EV-BC (ours)	101.33 \pm 14.37	38.3 \pm 6.81	138.8 \pm 18.6	22.0 \pm 6.1

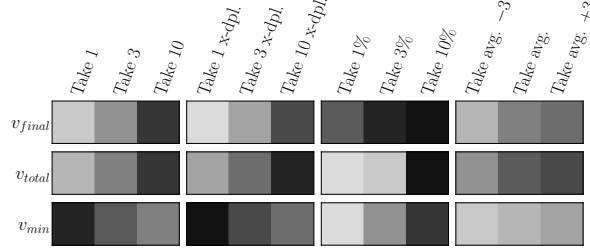


Figure 7: Colour-coded ordering of EVs for agents with varying behaviors in Tragedy of the Commons. The brighter, the higher an agent’s contribution to a given value function.

C TRAGEDY OF THE COMMONS EXPERIMENTS

We model ToC as a multi-agent environment where agents consume from a common pool of resources x_t , which grows at a fixed rate $g = 25\%$ at each time step t : $x_{t+1} = \max((1 + g) \cdot x_t - \sum_i c_{ti}, 0)$, with c_{ti} as the consumption of the i th agent at time t and $x_0 = 200$ as the starting pool. Hence, if all resources are consumed, none can regrow and no agents can consume more resources. The Tragedy of the Commons (ToC) environment features 4 different behavior patterns: *Take-X* consumes X units at every timestep, *Take-X-x-dpl* consumes X units if this does not deplete the pool of resources, *Take X%* consumes X% of the available resources, and *TakeAvg* consumes the average of the resources consumed by the other agents at the previous timestep (0 in the first timestep). For the small-scale experiment of 12 agents, we consider three agents for each pattern, with X values selected from the set 1, 3, 10. For the large-scale experiment of 120 agents, we simply replicate each agent configuration 10 times. We simulate both experiments for groups of size 3 and 10 respectively. We generate a simulated dataset using agents with four different behavior patterns. We first collect a dataset of observations for a small-scale experiment of 12 agents and simulate ToC for groups of three agents for 50 time steps (we later consider a group of 120 agents).

Due to the continuous nature of the state and action spaces in ToC, we first discretize both and then apply the same clustering methods used in the Overcooked scenario. We proceed by computing EVs for all agents as done in Overcooked (see Figure 3 for results). We implement imitation policies by replicating the averaged action distributions in the discretized states.

D COMPUTATIONAL DEMAND AND REPRODUCIBILITY

We used an Intel(R) Xeon(R) Silver 4116 CPU and an NVIDIA GeForce GTX 1080 Ti (only for training BC, EV2BC, group-BC, and OMAR policies). In Overcooked, generating a dataset took a maximum of three hours, and estimating EVs from a given dataset takes a few seconds. Behavior clustering consumes a couple of minutes, while Variance clustering took up to two hours per configuration (note that it is run 500 times). Training of the BC, group-BC, and EV2BC policies took no more than 30 minutes (using a GPU), while the OMAR baseline was trained for up to 2 hours. In Tragedy of Commons, each rollout only consumes a couple of seconds. Clustering times were comparable to those in Overcooked. Computing imitation policies is similarly only a matter of a few minutes.

As this submission is public, we will release the code for all experiments with the camera-ready version.