

An approximate embedding for designing ethical reinforcement learning environments

Arnau Mayoral-Macau^{a,*}, Manel Rodríguez-Soto^a, Enrico Marchesini^b, Martí Sánchez-Fibla^a,
Maite López-Sánchez^c, Juan Antonio Rodríguez-Aguilar^a and Alessandro Farinelli^d

^aIIIA-CSIC, Spain

^bMassachusetts Institute of Technology, Boston, USA.

^cUniversitat de Barcelona

^dDepartment of Computer Science, University of Verona, Verona, Italy

Abstract. This paper introduces the Approximate Ethical Embedding Process, an algorithm for automating the design of ethical environments for learning agents. Our algorithm helps build environments wherein multiple agents learn policies that align with an ethical (moral) value while simultaneously pursuing their individual objectives. Therefore, we contribute to endowing *environment designers* with algorithmic tools for building ethical environments. We demonstrate the ethical design process for two different settings of an environment where agents have to adhere to beneficence to promote the collective survival of the population. Our experiments show that our approximate embedding process successfully generates environments that incentivise the learning of value-aligned policies.

1 Introduction

As autonomous agents gain more prevalence in daily tasks [39, 13, 42, 5], their risks become more apparent. Thus, various international initiatives, such as the AI Act [8], mandate the systems to behave aligned with human values [28, 32, 9]. Thereby, as Multi-Agent Reinforcement Learning (MARL) algorithms have found application in diverse domains, they have also been used to instil *value-alignment* in the context of Machine Ethics [44, 27].

Machine Ethics pursues that (ethical) value-aligned behaviour involves proactivity in performing good (praiseworthy) actions (e.g., [9]). To achieve such an objective, the literature on Machine Ethics has extensively used Reinforcement Learning (RL) to help agents learn to behave ethically. Specifically, it is common in the literature [21, 1, 41, 17, 3, 34] to adopt an *agent-centered* approach to value alignment: an agent is guided towards a value-aligned behaviour by providing it with extrinsic, manually-tuned ethical rewards incorporated into its learning environment.

Alternatively to the agent-centred approach, we find in the literature recent *environment-centred* approaches to value alignment [23, 24]. This line of research takes an *environment designer* perspective, which focuses on automating the design of an ethical environment — the so-called ethical embedding process — for either a single agent [23] or multiple agents [24]. In such ethical environments, agents are *guaranteed* to learn ethical policies. Environment-centred approaches to value alignment are particularly appealing with respect to agent-centred approaches because they provide formal guarantees

regarding the learning of ethical policies and the automation of the reward design. However, these approaches are based on strict formal assumptions (e.g. full observability, convergence to optimality while learning), severely compromising their scalability.

Against this background, our goal is to contribute to the applicability of environment-centred approaches to value alignment. Thus, our main contribution is a new ethical embedding algorithm (henceforth *approximate embedding*) that goes beyond the scalability of the embedding algorithm in [24] (henceforth *optimal embedding*). There are major differences between our approximate embedding and the previous optimal embedding. First, our approximate embedding adopts more realistic assumptions (e.g., partial observability and no need for RL algorithms with convergence guarantees). Second, the approximate embedding builds on an entirely different technical method compared with the optimal embedding.

Like optimal embedding, our approximate embedding takes as input a multi-objective environment, with *ethical* and *individual* objectives (rewards), to produce a single-objective ethical environment wherein agents learn. This approach prevents agents from learning unethical policies by ensuring that it is in their best interest to behave ethically. Moreover, it implements two key computations differently from optimal embedding: (1) computing the *reference (ethical) joint policy* for agents to learn in an ethical environment; (2) computing the *ethical weight* to combine ethical and individual rewards into a single reward so that the ethical joint policy is the optimal policy to learn in the ethical environment.

First, our approximate embedding computes a reference ethical joint policy as a Nash equilibrium of a multi-agent multi-objective environment using Deep Reinforcement Learning (DRL). DRL has impressive results in approximating Nash equilibria despite its lack of theoretical guarantees [2, 43, 14]. Second, computing the ethical weight calls for introducing the Ethical Weight Finder (EWF), a novel algorithm based on binary search. More precisely, we make the following key contributions:

First, we present a novel embedding algorithm for building multi-agent ethical environments that incentivise the learning of ethical policies, the *Approximate Embedding*. Our algorithm builds upon two novel multi-objective RL developments:

- An algorithm, Multi-Agent LPPO (MALPPO), for computing the reference ethical policy. MALPPO is a multi-agent extension of the lexicographic proximal policy optimisation (LPPO) [31],

* Corresponding Author. Email: arnau.mayoral@iiia.csic.es

a state-of-the-art multi-objective RL algorithm. MALPPO computes a reference ethical joint policy as an (approximated) Nash equilibrium of an ethical MOPOMG.

- An algorithm, the Ethical Weight Finder (EWF), for computing the ethical weight to combine ethical and individual rewards into a single-reward, ethical environment. EWF searches for the ethical weight so that the optimal policy in the ethical environment is the reference ethical policy computed by MALPPO.

Then, we empirically show that approximate embedding successfully builds ethical environments for a large environment for which optimal embedding is inapplicable: an ethical version of the Gathering Game [15, 11] with more than 10^{64} states. In the resulting ethical environments, agents learn to adhere to the moral value of beneficence to ensure the survival of the whole agent population.

2 Background

The MARL literature formally defines a multi-agent environment as a *Markov Game* (MG) [2]. MGs are sequential decision-making settings where agents simultaneously act to modify the environment state and accumulate individual rewards. When agents have limited sensing capabilities over states, it is characterised as a *Partially Observable Markov Game* (POMG) [2]:

Definition 1 (POMG: Partially Observable Markov Game). A *Partially Observable Markov Game* is defined as a tuple $\langle S, A^{i=1,\dots,n}, R^{i=1,\dots,n}, T, O^{i=1,\dots,n}, \mathcal{O}^{i=1,\dots,n}, \gamma \rangle$. Here, S is a finite set of states, and each A^i represents the set of actions available to agent i . The transition function $T : S \times A^{i=1,\dots,n} \times S \rightarrow [0, 1]$ defines the probability of moving from state s to the next state s' , given the joint action $a = \langle a^1, \dots, a^n \rangle$ of all agents. For each agent i , the reward function $R^i : S \times A^{i=1,\dots,n} \times S \rightarrow \mathbb{R}$ specifies the reward r^i after applying joint action a to state s and transitioning to state s' . $O^{i=1,\dots,n}$ is a finite set of observations and the function $\mathcal{O}^{i=1,\dots,n} : A^{i=1,\dots,n} \times S \times O^{i=1,\dots,n} \rightarrow [0, 1]$ represents the probabilities over the agent's possible observations o^i given the state s and a joint action a . Finally, $\gamma \in (0, 1]$ is the discount factor which indicates how important future rewards are on the current state.

Each agent i aims to learn a policy (i.e., a behaviour) π^i that maximises its expected discounted accumulation of rewards $V^{\pi^i}(s) = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r^t]$, called the *value*¹ of the policy π^i at state s . Knowing such value function allows an agent to choose actions that lead to the most valuable states. However, computing an exact value function is challenging for environments with large state and action sets.

Typically, a joint policy $\pi = \langle \pi^1, \dots, \pi^n \rangle$ that maximises the return for all agents does not exist. Consequently, following the game theory literature, the main solution concept in MARL is reaching a *Nash Equilibrium* (NE) [2], defined as a joint policy in which no agent can unilaterally improve its current accumulation of rewards:

Definition 2 (Nash equilibrium). Given a *Partially Observable Markov Game* \mathcal{M} , a *Nash equilibrium* is a joint policy $\langle \pi_*^i, \pi_*^{-i} \rangle$ satisfying that for every agent i and every state observation $\vec{o} = (o^1, \dots, o^n)$, with each o^i in O^i , the policy π_*^i of agent i is a *best-response* against $\pi_*^{-i}(s)$, that is, it maximises the return against the joint policy of the rest of the agents π_*^{-i} :

$$V_{\langle \pi_*^i, \pi_*^{-i} \rangle}(\vec{o}) \geq V_{\langle \pi^i, \pi_*^{-i} \rangle}(\vec{o}), \text{ for every } \pi^i \text{ and } \forall \vec{o} \in \vec{O}, \quad (1)$$

¹ Here, “value” refers to a metric used in RL to evaluate the utility of states or actions. This term is not related to moral or ethical values.

where $V_{\pi^i}^i(\vec{o})$ is the expected discounted accumulation of rewards $\mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r^t \mid \pi, \vec{o}]$ of agent i if all agents follow the joint policy $\pi = \langle \pi^1, \pi^{-i} \rangle$ after observing \vec{o} .

Computing Nash equilibria in Markov games is a complex but well-studied problem [18, 40, 4]. The exact computation of equilibria is computationally intensive, with complexity increasing rapidly in high-dimensional environments. Recently, state-of-the-art single-agent algorithms leveraging deep neural networks—extended to multi-agent scenarios—have been employed to approximate equilibria under paradigms such as independent learning and centralised training with decentralised execution [43, 20, 6].

Multi-Objective MARL. When considering multiple ($m > 1$) learning goals in a POMG, agents aim to maximise the accumulation of rewards obtained with respect to each of the objectives in the so-called *Partially Observable Multi-Objective Markov Game* [29]:

Definition 3 (MOPOMG: Multi-Objective Partially Observable Markov game). A (finite) *Partially Observable m -objective Markov game* of n agents is defined as a tuple $\mathcal{M} = \langle S, A^{i=1,\dots,n}, \vec{R}^{i=1,\dots,n}, T, O^{i=1,\dots,n}, \mathcal{O}^{i=1,\dots,n}, \gamma \rangle$ whose elements $S, A^{i=1,\dots,n}, T, O^i, \mathcal{O}^i$, and γ are defined exactly like those of an POMG. On a MOPOMG, the reward function $\vec{R}^{i=1,\dots,n} = (R_1^i, \dots, R_m^i)$ is vectorial, where each scalar reward function $R_j^i \in \vec{R}^{i=1,\dots,n}$ is the reward function regarding the j -th objective. Accordingly, at each simulation time step, the agent i gets a vectorial reward signal $\vec{r}^i = (r_1^i, \dots, r_m^i)$.

In a MOPOMG \mathcal{M} , the value \vec{V}^π of a policy π is a vector. Comparing vectorial returns requires additional information about objective priorities. This information, such as a user's prioritisation among objectives, is crucial for determining an optimal solution. A weight vector $\vec{w} \in \mathbb{R}^m$ can represent this prioritisation, allowing the scalarisation of value vectors as $\vec{w} \cdot \vec{V}^\pi$ into comparable scalar values. The single-agent MORL literature [26] defines the *convex hull* (CH) as the set of policies that are optimal for some scalarisation weight vector \vec{w} . Therefore, in the single-agent literature, the CH is a solution set containing the optimal policies for any linear prioritisations a final user might prefer. However, the CH has been extended to multi-agent scenarios *only* when agents fully cooperate and share exactly the same objective [25].

Besides linear scalarisation with a weight vector, the multi-objective literature also considers non-linear prioritisations. Lexicographic orders (LO) [35] are explicit orderings of the objectives where objectives are prioritised over those that follow them on the ordering. Thus, the two policies are compared objective-wise, following the LO. Imagine a policy π_1 that receives a return $\vec{V}^{\pi_1} = (3, 4, 5)$ and π_2 with $\vec{V}^{\pi_2} = (3, 5, 3)$. Then, for an LO $\ell_1 = \{V_1 \succeq V_2 \succeq V_3\}$, policy π_2 is better, as it has equal V_1 and has greater V_2 ; under another LO $\ell_2 = \{V_3 \succeq V_1 \succeq V_2\}$, π_1 is optimal over π_2 as it gets more return for V_3 .

3 The Embedding Problem

While the literature on automated environment design has primarily focused on value alignment, the multi-objective embedding problem can be formalised for any pair of objectives. This approach seeks to design environments in which agents learn to accomplish their primary tasks while simultaneously aligning with an additional *alignment* objective. This secondary objective may be ethical, safety-related, or any other objective that agents are expected to abide

by, in the same vein as [36, 37]. Crucially, the goal of the embedding problem is to design an environment where the *only* optimal behaviour is the intended, aligned behaviour. Hence, our problem amounts to creating environments where agents are incentivised to learn behaviours aligned with an extrinsic alignment objective. However, rather than developing a learning algorithm, we focus on transforming the agents’ learning environment by embedding the primary objective with the alignment objective. With this embedded reward function, agents do not have the means to separately prioritise the objectives when learning. This environment-centred perspective is motivated by the fact that such learning environments may be used by third-party entities, whose choice of learning algorithms and prioritisation of the objectives is unknown.

Despite the generality of the environment-centred approach, in this paper, we focus on the ethical design of environments. To design an *ethical* environment, our algorithm takes as input an initial (henceforth *source*) environment with two objectives: an individual task representing the agents’ primary goal and an *ethical objective* that evaluates and quantifies the agents’ alignment with a specific moral value. A moral value or ethical principle, in Ethics, represents a moral goal worth pursuing [38]. The environment is then transformed into a single-objective *ethical environment*, where the two objectives are embedded together in one reward function that prioritises the ethical objective. Therefore, maximising this single reward can only lead to value-aligned behaviour. Henceforth, we formalise the Ethical Embedding problem by defining: first, the agents’ source and ethical environments; and then the problem itself of how to transform an original environment into an ethical environment wherein agents learn to behave ethically.

We model the agents’ source environment as a two-objective *Partially Observable Markov Game* with: an individual reward function R_0^i (the reward function that rewards each agent i for fulfilling its individual objective), and an ethical reward function R_e^i that rewards each agent i when behaving ethically. Thus, to properly incentivise ethical behaviour, R_e^i has to be carefully designed following the ethics literature. Ethical frameworks that assess actions based on their consequences are particularly suitable when designing ethical reward functions like R_e^i . In our work, we employ the ethical framework presented in [22, 23, 24, 30] to construct ethical reward functions R_e^i , which is grounded in the Ethics literature. However, any other ethical framework for properly designing R_e^i could be used. We refer to such family of Markov Games as *Ethical Multi-Objective Partially Observable Markov Games*:

Definition 4 (Ethical MOPOMG). *An Ethical Multi-Objective Partially Observable Markov Game (EMOPOMG) \mathcal{M} is defined as a tuple $\langle S, A^{i=1,\dots,n}, R_0^{i=1,\dots,n}, R_N^{i=1,\dots,n}, R_E^{i=1,\dots,n}, T, O^{i=1,\dots,n}, \mathcal{O}^{i=1,\dots,n}, \gamma \rangle$ such that for each agent i : First, R_0^i is the individual reward function of each agent i , representing their individual objective. Then, $R_N^i : S \times \mathcal{A}^i \rightarrow \mathbb{R}^-$ is the normative reward function of each agent i , penalising blameworthy actions, i.e. committing morally prohibited actions or ignoring moral obligations. Finally, $R_E^i : S \times \mathcal{A}^i \rightarrow \mathbb{R}^+$ is the evaluative reward function of each agent i , rewarding praiseworthy or supererogative actions².*

Tuple elements $S, A^{i=1,\dots,n}, T, O^{i=1,\dots,n}, \mathcal{O}^{i=1,\dots,n}$, and γ of \mathcal{M} are defined identically to Partially Observable Markov Games.

Note that this framework can represent any ethical value that can

² Morally good but not mandatory actions
(<https://plato.stanford.edu/entries/supererogation/>)

be encoded as a set of moral prohibitions, obligations and recommendations over the agents’ actions.

We define an *ethical equilibrium* of an Ethical MOPOMG as an NE with respect to the ethical reward function $\vec{R}_e = \vec{R}_N + \vec{R}_E$, where \vec{R}_e denotes $R_e^{i=1,\dots,n}$. Among ethical equilibria, we highlight *best-ethical equilibria*. A best-ethical equilibrium π_* is a Nash equilibrium with respect to the individual reward function of agents, subject to also being an ethical equilibrium. Notice that the notion of Ethical MOPOMG and ethical equilibrium of an Ethical MOPOMG generalise analogous concepts in [24].

Thus, the *Ethical Embedding* problem is: how to design, from a given Ethical MOPOMG \mathcal{M} , a (single-objective) POMG \mathcal{M}_e that provides enough incentives to the agents to learn to behave ethically (a best-ethical equilibrium). These incentives are provided to agents by weighting ethical rewards with a large enough ethical weight $w_e > 0$ such that the (scalar) reward that each agent receives $R_0^i + w_e \cdot R_e^i$ promotes the agents to behave ethically. Formally, our *target* environment is an *Ethical Partially Observable Markov Game*:

Definition 5 (Ethical Partially Observable Markov Game). *Let \mathcal{M} be an Ethical MOPOMG with reward functions R_0^i, R_e^i for each agent i . We refer to the Ethical Partially Observable Markov Game \mathcal{M}_e associated with \mathcal{M} to a (single-objective) POMG with reward function $R_0^i + w_e \cdot R_e^i$ with $w_e > 0$, such that at least one Nash equilibrium of \mathcal{M}_e is a best-ethical equilibrium in \mathcal{M} .*

To ease notation, we mark $\mathcal{M}_{\langle w \rangle}$ to denote a single-objective POMG with rewards $R_0^i + w \cdot R_e^i$ resulting of scalarising \mathcal{M} ’s rewards (R_0^i, R_e^i) , with ethical weight w , for each agent i . Moreover, henceforth, we refer to an Ethical MOPOMG as a *source* environment and to its respective Ethical POMG as the *target* environment.

4 Solving the Embedding Problem

Figure 1 outlines the approximate embedding process. Following the environment-centred perspective previously mentioned, this process aims to design an ethical environment by transforming a multi-objective environment (a source environment \mathcal{M}) into an *ethical single-objective* environment (the target environment \mathcal{M}_e). Our purpose is that in the ethical environment \mathcal{M}_e , any third-party agent, independently of its learning algorithm, will learn a value-aligned behaviour. The approximate embedding process follows three steps:

1. **Reference policy computation.** We compute a so-called *reference joint policy* π_r in the source environment. This is the ethical policy we want agents to, ultimately, jointly learn in the ethical single-objective target environment.
2. **Ethical weight computation.** We compute the ethical weight w_e to transform the source environment into our target environment.
3. **Ethical environment synthesis.** We build the target environment as a POMG ($\mathcal{M}_e = \mathcal{M}_{\langle w_e \rangle}$) by scalarising the ethical rewards in the original environment using the ethical weight w_e computed at the previous step.

4.1 Reference policy computation

The first thing we must compute to obtain our target environment is an ethical policy that serves as a reference. Using the value vectors (multi-objective returns) of this policy, we will be able to properly design an ethical environment where ethical behaviour is preferred over individual behaviours. This reference policy is formally a best-ethical equilibrium in the source environment, the MOPOMG. As defined in Section 3, this best-ethical equilibrium corresponds to the

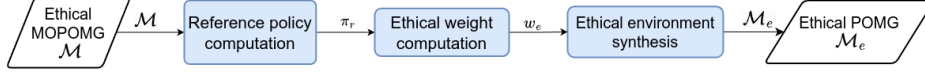


Figure 1. Approximate Ethical Embedding Process.

equilibrium agents reach when prioritising ethical rewards over individual rewards. This reference policy π_r is important for two reasons. First, it is the joint policy we want to incentivise the agents to learn in the target environment. Second, we will need it to find a scalarisation weight, the so-called ethical weight w_e , to combine ethical and individual rewards in the target environment.

To obtain such an ethical reference policy, we need a learning algorithm that always prioritises ethical returns over individual returns. In the MORL literature, this type of non-linear prioritisation is studied under the name of lexicographic RL (LRL) [10, 31, 35]. LRL prioritises objectives according to an explicit ordering.

Using a lexicographic learning algorithm in an ethical MOPOMG \mathcal{M} , which prioritises the ethical objective over the individual objective ($R_e \succeq R_0$), ensures the learned joint policy first optimises the ethical objective and then, without altering their obtained ethical return, optimise the individual objective. Notice that, when learning a joint policy in an ethical MOPOMG with a lexicographic algorithm, any Nash equilibrium will be a *best-ethical* equilibrium by definition.

Consequently, by using a lexicographic learning algorithm in the source environment \mathcal{M} , we can obtain ethical joint policies that abide by the ethical objective. This reference policy, π_r , represents the behaviour we aim to establish as the optimal policy in the resulting scalarised ethical environment \mathcal{M}_e , so that any agent learning in \mathcal{M}_e finds it optimal to align with the moral value encoded in R_e .

However, no such algorithm exists for multi-agent environments. Subsequently, to learn a best-ethical equilibrium, we implemented a new lexicographic multi-agent algorithm. We provide more details on this contribution in Sections 5 and 6.2. The next section shows how to exploit a reference policy π_r to compute the ethical weight.

4.2 Ethical weight computation

Recall that our goal is to build a target environment \mathcal{M}_e whose optimal policy for any learning algorithm is the (ethical) reference policy π_r . Since the target environment \mathcal{M}_e is intended to be single-objective, we must find an ethical weight w_e that scalarises the ethical reward in the original environment \mathcal{M} and sets π_r as the equilibrium in the target environment \mathcal{M}_e .

We can exploit the multi-objective returns of the reference policy \vec{V}^{π_r} to compute the ethical weight. Since we know that the reference policy is the one that attains the largest ethical return, the ethical weight w_e that we choose to scalarise the source environment \mathcal{M} must be large enough to make π_r the optimal policy in the scalarised environment \mathcal{M}_e . More technically, we must compute an ethical weight w_e such that the Nash equilibrium π_r attains more scalarised return than any other equilibrium in \mathcal{M}_e . Thereafter, we will be able to build our target environment \mathcal{M}_e using w_e wherein the (ethical) reference policy stands as the optimal policy to learn.

When computing the ethical weight w_e , we target an ethical weight w_e as low as possible while still capable of incentivising the learning of ethical equilibria. There are two reasons for searching for such an ethical weight. First, a large ethical weight could hinder or even prevent agents from learning their individual objectives, as demonstrated in the experiments section. Second, we consider that a reward function might have an associated cost when deploying the

agents. Thus, an excessive weight w_e would involve a higher cost.

Since we know that an ethical equilibrium π_r prioritising ethical reward exists, we assume there must be an ethical weight w_e sufficiently large such that when learning with scalarised objectives $R = R_0 + w_e \cdot R_e$, the ethical objective is completely prioritised over the individual objective. With that assumption in mind, we argue that, in the space of possible weights, we can differentiate two intervals: (i) $[0, w_e)$ with weights that do not incentivise the learning of an ethical equilibrium enough, and (ii) $[w_e, \infty)$ with ethical weights that effectively incentivise ethical behaviour. Note that any weight w in the interval $[w_e, \infty)$ allows building a scalarised, target environment whose ethical equilibrium is the reference policy. However, for the reasons mentioned above, we pursue the lowest possible ethical weight w_e .

To identify w_e , we need a method for testing whether a given weight w qualifies as an ethical weight. Once a reference policy π_r is computed, testing whether a weight w is ethical is straightforward: (1) build a scalarised environment $\mathcal{M}_{(w)}$; (2) learn the optimal joint policy π in $\mathcal{M}_{(w)}$ (the so-called *approximate reference policy*); (3) compare the ethical returns of π and π_r . If π achieves *close enough* ethical return to the reference policy’s ethical return (i.e., $|V_e^\pi - V_e^{\pi_r}| < \tau$ for some *policy approximation error* $\tau > 0$), then we consider w as an ethical weight. We do so because it sufficiently incentivises the learning of a policy that is as ethical as the reference policy. For simplicity, henceforth, we will refer to testing whether a weight w is ethical as function $test(\mathcal{M}_{(w)}, \tau)$.

This definition for a single-agent $test(\mathcal{M}_{(w)}, \tau)$ function can be extended to multi-agent setups simply by adding a sum over the agents returns such that $\left| \sum_{i=1}^n [V_e^{\pi_i} - V_e^{\pi_r^i}] \right| < \tau$.

By repeatedly using the *test* function, we get valuable information on the value of w_e needed to create an environment whose equilibrium is π_r . Therefore, the set of joint policies obtained after testing different weights on a MOPOMG is relevant to inform our search for the minimum ethical weight. To characterise this set more rigorously, we introduce the *Nash Convex Hull (NCH)* for a MOPOMG, which we define as the set of joint policies that include a Nash equilibrium for any possible linear scalarisation weight vector. Formally:

Definition 6 (Nash Convex Hull). *The Nash Convex Hull of a MOPOMG \mathcal{M} is the set of policies that contains the Nash equilibrium of all possible scalarisations $\mathcal{M}_{(\vec{w})}$ performed with any linear scalarisation weights \vec{w} : $NCH(\mathcal{M}) = \{\pi \in \Pi^{\mathcal{M}} | \exists \vec{w} : \pi \in NE(\mathcal{M}_{(\vec{w})})\}$, where $\Pi^{\mathcal{M}}$ is the set of possible policies in \mathcal{M} and $NE(\mathcal{M}_{(\vec{w})})$ is the set of equilibria in the scalarised POMG $\mathcal{M}_{(\vec{w})}$.*

To obtain the ethical weight, there is no need to compute the whole *NCH*. Consequently, the resulting set of our search is a subset of the *NCH* that contains the Nash equilibria associated with the candidate ethical weights considered throughout the search for the minimal ethical weight. Additionally, considering that the learning algorithm might lack guarantees to find exact equilibria, we frame the set of policies we compute as an *approximate NCH*. Later in this section, we detail the construction of this *NCH*.

With $test(\mathcal{M}_{(w)}, \tau)$ and considering a finite search space between 0 and an upper bound large enough, we can use any search algorithm

to find our desired ethical weight. We propose using a general search paradigm like binary search [16] to automate the search of an *ethical weight* inside a search interval $w_e \in I$. Then, our weight computation would work as follows. Let $I = [w_l, w_r]$ represent the unexplored search space, where w_l is not an ethical weight and w_r is an ethical weight. We can set w_l to 0 because it produces an environment without ethical rewards, and hence no ethical policies. Regarding w_r , the environment designer can set it to any large upper-bound number based on their expert knowledge. Thus, the search would begin with an initial, non-negative solution $w_e = w_r/2$, along with a precision parameter ϵ , which controls the depth of the search, and a policy approximation parameter τ . We test with $\text{Test}(\mathcal{M}_{(w_r/2)}, \tau)$ whether $w_r/2$ is an ethical weight. That is, we test whether the optimal policy in $\mathcal{M}_{(w_r/2)}$ is as ethical as the reference policy. If it is, we update the interval to $I = [w_l, w_r/2]$. If not, we update it to $I = [w_r/2, w_r]$. The procedure continues iteratively narrowing interval I until the distance between its endpoints is less than or equal to ϵ , while ensuring that the left endpoint remains a non-ethical weight and the right endpoint remains an ethical weight. The right endpoint is the approximation to the minimum ethical weight we are looking for. We refer to the whole process as the *Ethical Weight Finder* (EWF). It is important to note that this algorithm does not compute the exact w_e that makes π_r the optimal policy. Instead, it constructs an approximate subset NCH around such theoretical w_e . Since we progressively narrow the interval containing the exact w_e to a width of ϵ , by the end of the search, we obtain a set of policies that are Nash equilibria for weights close to w_e . As environment designers, we then select the weight that best incentivises ethical behaviour.

So far, we have demonstrated how to perform a binary search to identify the desired ethical weight w_e , minimising the number of times we run a learning algorithm, as this is the most computationally expensive part of the search. The approach described so far is only informed by the value vectors. We tried to develop a method as generic as possible, following the lines of [33]. However, we admit that the search can be optimised in different ways that may result in a more optimal search for specific environments.

5 Implementing approximate embedding

Next, we detail how to compute the two main elements of our approximate embedding, the reference joint policy and the approximate reference policies, to find our desired ethical weight.

For the sake of understanding, we begin with the computation of the approximate reference policies, computed at each iteration of our EWF. Recall that the approximate reference policies are formally Nash equilibria in a *single-objective* POMG. Thereafter, as explained in Section 4.2, we must compute several Nash equilibria in multiple scalarised POMGs during the search for the ethical weight. We can find them using single-objective Multi-Agent Deep Reinforcement Learning (MADRL) techniques. Indeed, MADRL techniques have shown great results in approximating equilibria for complex cooperative games [7, 43, 6] by extending powerful single-agent algorithms to MARL scenarios. Examples include independent PPO (IPPO), and its centralised training and decentralised execution (CTDE) counterpart, MAPPO [2]. Here we choose to employ MAPPO because its parameter sharing and use of centralised critic help reduce sample complexity, increase learning efficiency, and stabilise learning.

On the other hand, the reference policy corresponds to a best-ethical equilibrium in a *multi-objective* POMG. As such, its computation requires a multi-objective RL algorithm. Crucially, lexicographic algorithms can be used to compute a best-ethical equilibrium when priority is given to the ethical objective. Recently, [31]

introduced both value-based and policy-based lexicographic RL algorithms, including LPPO. Nonetheless, those algorithms only exist and were tested for the single-agent case.

Due to the lack of multi-agent lexicographic algorithms, we have developed our so-called MALPPO, a novel CTDE extension of LPPO. MALPPO adopts the CTDE paradigm from MAPPO, which features a centralised critic with access to all agents' observations. The critic's estimations are used by each agent to independently train its own policy using LPPO. We implemented our MALPPO in a fork of the EPyMARL library [19], extended for multi-objective algorithms³.

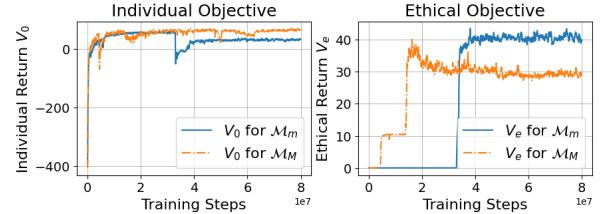


Figure 2. Sum of the individual and ethical returns of all agents for the learned reference policies in environments \mathcal{M}_m (solid line) and \mathcal{M}_M (dash-dot line).

6 Empirical analysis

In this section, we employ our approximate embedding to build ethical environments for the *Ethical Gathering Game* (EGG), an ethical extension of the Gathering Game originally introduced in [15] and outlined in Section 6.1. Importantly, we empirically verify that the ethical environments that we design do incentivise the learning of ethical joint policies in large environments (Section 6.3). Moreover, we quantify the *price to pay*, the cost, to fabricate ethical environments through approximate embedding (Section 6.2).

6.1 The Ethical Gathering Game

The *Ethical Gathering Game* is a grid-world environment where agents with different capabilities gather apples to survive [24]. An agent survives if it manages to accumulate apples beyond some survival threshold. The EGG is an ethical extension of the well-known Gathering Game [15] in the MARL literature. The EGG simulates an unequal environment where only efficient agents can survive by themselves while inefficient agents require the efficient agents' help.⁴ The EGG also introduces a *public* donation box: agents can either donate apples to the box or retrieve apples from it. Then, two objectives drive agents in the EGG: an *individual* objective to collect apples for personal survival, and an *ethical* objective to contribute to collective survival. As environment designers, we will create ethical environments whose agents learn the moral value of *beneficence* to help the whole agent population survive. For that, we will instill alignment with beneficence through approximate embedding.

Source environments. We built two different source environments: one with a majority of inefficient agents (\mathcal{M}_m), and another one with a minority of inefficient agents (\mathcal{M}_M). The settings for each environment are such that inefficient agents cannot survive unless the efficient agents learn to behave ethically to help them. Thus, the survival of the whole population is possible in both environments, though, as we will show, the amount of required beneficence varies. This difference is key, as it makes the environments totally different with respect to computing the embedding, as the minimum ethical weight

³ <https://github.com/maymac00/mo-epymarl.git>

⁴ The Gathering Game [15] deals with resource depletion, a different problem from ours, for which agents must coordinate to survive.

we look for will be different. More precisely, each environment has 5 agents and a 15-apple donation box, and the individual survival threshold is 40 apples. \mathcal{M}_m has 40% (minority) of efficient agents and \mathcal{M}_M has 80% (majority) of efficient agents. Importantly, since the EGG uses the same grid size as the original Gathering game, our settings configure two different large environments with more than 10^{64} states each. Therefore, optimal embedding [24], which relies on Q-Learning, cannot be applied to these environments.

MOPOMG formalisation. The EGG is a MOPOMG whose agents act during a 500-step gathering season. To specify a *source environment*, following Def. 4 of Ethical MOPOMG, we define an individual reward function R_0^i that gives a reward of +1 to an agent for each collected apple. There are two types of agents: an inefficient agent has a 15 % chance of gathering an apple from the ground, while an efficient agent has an 85 % chance. An agent receives a negative reward $R_0^i = -1$ when: (i) being below the survival threshold at a given time step, or (ii) donating an apple to the public donation box. As to ethical reward, $R_e^i = R_N^i + R_E^i$, we penalise *unethical actions* ($R_N^i = -1$ when taking an apple from the donation box despite having enough to survive), and positively reward *ethical actions* ($R_E^i = 0.7$ when donating an apple when having enough to survive, namely accumulated apples beyond the survival threshold).

Metrics. Once designed an ethical environment, we have the agents independently learn their policies. We will then employ two metrics to analyse whether they have learned to behave ethically: (i) the **survival rate** (how many times *all* agents are able to reach the survival threshold on average); (ii) the **ethical returns** (the accumulation of discounted ethical rewards that all agents obtain on average).

6.2 Designing ethical environments

Approximate embedding algorithms. As discussed in Section 5, we will employ MALPPO to compute reference policies and standard MAPPO within the EWF algorithm to find the ethical weight. We run MALPPO and MAPPO using the EPyMARL library [19].

Learning reference policies. We applied approximate embedding to our two source environments: *minority* (\mathcal{M}_m) and *majority* (\mathcal{M}_M). We first computed the reference joint policy for each environment using MALPPO, our MARL lexicographic algorithm. Figure 2 presents the expected (individual and ethical) returns that the reference policy obtained in each environment. We summed up the returns over all agents and smoothed them using an exponential moving average at 0.6 to reduce variance visually.

Finding ethical weights. We applied the EWF algorithm to search for the ethical weight for our two EGG environments. Recall that an ethical weight allows us to scalarise ethical rewards from a source environment so that the reference (ethical) policy becomes optimal in an ethical, target environment. For both environments, we used $[0, 10]$ ($w_l = 0, w_r = 10$) to search for ethical weights, $\tau = 4.0$ as policy approximation error, and $\epsilon = 0.2$. Now, for each environment, we search for an as-low-as-possible ethical weight w_e that produces a scalarised environment $\mathcal{M}_{\langle w_e \rangle}$ whose Nash equilibrium is ethical, namely a good approximation of the reference policy. Figure 3 shows the evolution of the search for the ethical weight in each environment for six iterations. For iteration i and ethical weight w_e (different for each i), we show the multi-objective expected returns of policy i (blue circle), which we previously computed with MAPPO as a Nash equilibrium in the scalarised environment $\mathcal{M}_{\langle w_e \rangle}$. As iterations increase, w_e evolves as well as the (approximate) Nash equilibrium in each $\mathcal{M}_{\langle w_e \rangle}$. The Nash equilibrium for each w_e can be regarded as an approximation to the reference policy, and the set of policies that we obtain, as an approximate *NCH* as defined in Def. 6.

Notice that the *approximate* reference policies within each shaded grey area can be considered ethical because they are close enough, according to policy approximation error τ , to the reference policy ($|\sum_{i=1}^n [V_e^{\pi^i} - V_e^{\pi_r}]| < \tau$ considering all n agents). From the new set of policies, we finally select $w_e = 2.5$ and $w_e = 1.71875$ as ethical weights for \mathcal{M}_m for \mathcal{M}_M respectively because they are the lowest weights leading to good enough approximations of their reference policies. In detail, our EWF converged to these points because we obtained the final intervals $I_m = [2.34375, 2.5]$ for \mathcal{M}_m , and $I_M = [1.56, 1.71875]$ for \mathcal{M}_M , with both intervals with a length smaller than $\epsilon = 0.2$, and with the left extrema being unethical, and the right extrema being ethical. Notably, in Figure 3, we observe that using a weight larger than necessary (i.e., with $w_e \geq 5$) is detrimental to individual return while providing only a marginal improvement in ethical return. This underscores the importance of identifying the approximate minimal weight as we motivated in Section 4.2.

Embedding cost. To measure the cost of the design process, we analyse the number of algorithm executions required to compute: (i) the reference policy; and (ii) the approximate reference policies when searching for the ethical weights. For each source environment, we run MALPPO once (to compute each reference policy), and we run MAPPO six times (during the ethical weight search). The total amount of steps needed to learn the necessary policies to design ethical environments is 420M for both \mathcal{M}_m and \mathcal{M}_M . MALPPO needed 80M steps to learn both policies, while the different runs of MAPPO during both binary searches needed 70M on average.

We conclude that there is a price to pay to design an ethical environment. In this particular experiment, designing each EGG ethical environment required running our MARL algorithms seven times. This investment paid off because agents learned ethical policies in the EGG ethical environments, as we show next.

Minority environment \mathcal{M}_m				
Ag.	Ethical Return π_r	Ethical Return π_*	Individual Return π_r	Individual Return π_*
1	17.54 \pm 8.67	18.35 \pm 5.83	26.15 \pm 16.97	27.18 \pm 9.53
2	17.61 \pm 8.53	17.51 \pm 5.79	26.4 \pm 16.78	26.17 \pm 9.68
3	1.69 \pm 1.42	0.77 \pm 1.09	-8.01 \pm 10.87	-8.85 \pm 8.24
4	1.7 \pm 1.31	0.54 \pm 0.91	-7.83 \pm 10.64	-10.16 \pm 8.70
5	1.64 \pm 1.34	0.60 \pm 0.85	-7.63 \pm 10.52	-10.36 \pm 8.62
Eff.	17.575 \pm 0.035	17.93 \pm 0.42	26.27 \pm 0.13	26.67 \pm 0.50
Ineff.	1.676 \pm 0.026	0.63 \pm 0.09	-7.82 \pm 0.15	-9.79 \pm 0.67
Majority environment \mathcal{M}_M				
1	7.13 \pm 3.97	7.26 \pm 4.59	19.11 \pm 9.49	19.70 \pm 10.64
2	6.96 \pm 4.13	7.09 \pm 4.81	18.59 \pm 9.63	18.10 \pm 10.18
3	7.08 \pm 3.89	7.49 \pm 4.74	18.75 \pm 10.01	20.30 \pm 10.80
4	7.05 \pm 3.96	7.41 \pm 4.62	18.65 \pm 9.77	20.40 \pm 10.29
5	0.83 \pm 0.66	0.00 \pm 0.02	-7.93 \pm 5.86	-18.10 \pm 8.44
Eff.	7.05 \pm 0.06	7.34 \pm 0.26	18.775 \pm 0.2	19.92 \pm 1.16
Ineff.	0.83 \pm 0.00	0.00 \pm 0.00	-7.93 \pm 0.00	-18.10 \pm 0.00

Table 1. Value vectors of each agent for the reference policy π_r and the ethical policy π_* learned in the designed environments \mathcal{M}_m and \mathcal{M}_M . Results correspond to an average of 2000 rollouts. Recall that π_r is the same policy for all agents, but it is not the case for π_* .

6.3 Learning in the ethical environments

We built the ethical environments for the source environments \mathcal{M}_m and \mathcal{M}_M using the weights computed above. To prove that the environments are indeed ethical, we will let agents learn with a single-objective learning algorithm to test whether they learn to behave ethically. Specifically, we employed IPPO [7], where all agents learn independently within each environment. This approach allows us to demonstrate that in an ethical environment designed with approximate embedding, ethical behaviour emerges as the optimal policy for

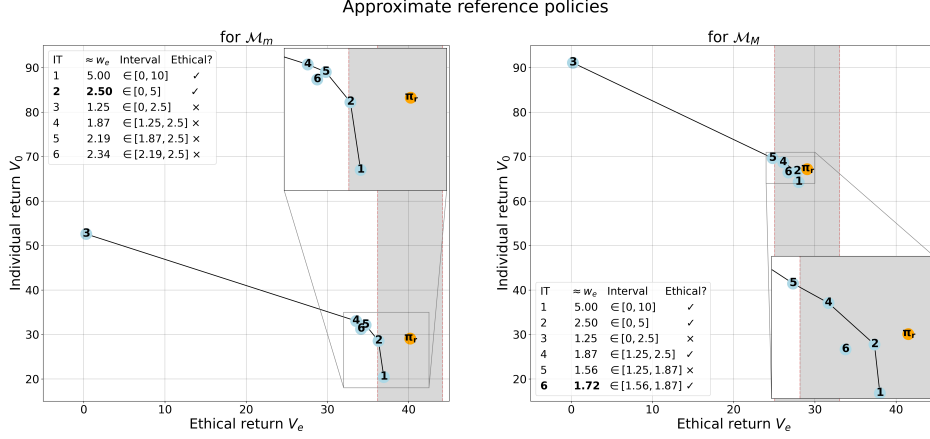


Figure 3. Approximate reference policies obtained in the search for the ethical weights of environments \mathcal{M}_m (left) and \mathcal{M}_M (right). Black lines show the approximate NCH around the exact ethical weight needed to build an ethical environment.

all agents, even in the absence of an explicit coordination mechanism during the training phase, such as the ones provided by MAPPO.

Policy	Survival rate in \mathcal{M}_m	Survival rate in \mathcal{M}_M
π_r	100%	99%
π_*	100%	98%
π_u	10%	66%

Table 2. EGG metrics for the reference (π_r), ethical (π_*), and unethical (π_u) policies. The expected ethical return of the policies is shown only for the efficient agents of the population.

The results obtained in the designed environments show that, indeed, agents were incentivised to learn a joint policy as ethical as the reference policy, as the metrics in Table 2 illustrate. First of all, the collective survival is achieved almost 100% of the times for both the reference policy π_r and the joint policy π_* that agents learn in both target ethical environments (column 3 of Table 2). For comparison, an unethical policy π_u , trained in an environment with ethical weight $w_e = 0$, reaches a much lower level of collective survival. Moreover, in terms of ethical returns, joint policies learnt in the target ethical environments obtain very close average returns to those obtained by the reference policies (column 4 of Table 2).

Considering all agents’ returns, not only efficient agents, Table 1 shows the value vectors of the ethical policy π_* and the reference policy π_r for all agents. Overall, we can see how the value vectors are within the tolerance parameter $\tau = 4$. Interestingly, the policies π_* learned in the designed environments achieve better individual returns than the reference policy. When looking at the statistics of the ethical gathering game rollouts, we find that efficient agents in π_* learned to gather more apples than those in π_r in the same amount of time, which increased their return and led to inefficient agents not gathering as much as they did with π_r .

Note that both policies have been achieved through different algorithms. While in π_r , agents shared the same value and policy network, in π_* , each agent learned its own value and policy networks. This can explain the difference between the deviations for agents sharing efficiency groups in both policies.

Finally, Figure 4 helps us understand the policies that inefficient and efficient agents learn in the majority, \mathcal{M}_M , environment. Agents learn analogous policies in the minority, \mathcal{M}_m , environment. The figure shows the median number of apples (over 2000 policy rollouts) that efficient and inefficient agents have throughout the episode. We observe that, in all runs, efficient agents have learned to donate their

surplus apples to aid inefficient agents in survival. Moreover, agents collect apples from the donation box only when they do not have enough apples to survive. This behaviour confirms again that, indeed, agents have been incentivised to learn to behave ethically.

In summary, all results from the experiments we conducted corroborate that the approximate embedding algorithm can design environments where agents learn best-ethical policies.

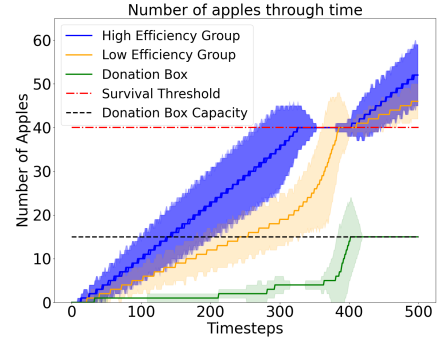


Figure 4. Median number of apples (and inter-quartile ranges) collected by agents throughout 2000 episodes of 500 steps in \mathcal{M}_M .

7 Conclusions and Future Work

This work introduced the *Approximate Embedding*, an algorithm for designing environments where all agents are incentivised to learn to behave ethically. Our empirical analysis shows that by combining deep RL and novel MORL tools like MALPPO, our *approximate embedding* successfully incentivises the learning of ethical policies in a large environment such as the *Ethical Gathering Game*. As future work, we aim to reduce the computational costs of approximate embedding and to evaluate it in further environments. Nevertheless, this will first require the engineering of more MARL environments with ethical objectives because they do not currently exist. For instance, we will consider an ethical reformulation of the cleaning game [12].

8 Acknowledgements

The research presented in this paper was supported by the EU-funded VALAWAI (# 101070930) project, and the Spanish-funded VAE (# TED2021-131295B-C31), EMOROB CARE (# IASOMMA2024), and ACISUD (PID2022-136787NB-I00) projects.

References

- [1] D. Abel, J. MacGlashan, and M. L. Littman. Reinforcement learning as a framework for ethical decision making. In *AAAI Work.: AI, Ethics, and Society*, volume 92, 2016.
- [2] S. V. Albrecht, F. Christianos, and L. Schäfer. *Multi-Agent Reinforcement Learning: Foundations and Modern Approaches*. MIT Press, 2024. URL <https://www.marl-book.com>.
- [3] A. Balakrishnan, D. Bouneffouf, N. Mattei, and F. Rossi. Incorporating behavioral constraints in online ai systems. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33:3–11, 07 2019. doi: 10.1609/aaai.v33i01.33013.
- [4] L. Buşoniu, R. Babuška, and B. De Schutter. Multi-agent reinforcement learning: An overview. *Innovations in multi-agent systems and applications-1*, pages 183–221, 2010.
- [5] Y. Chang, X. Wang, J. Wang, Y. Wu, K. Zhu, H. Chen, L. Yang, X. Yi, C. Wang, Y. Wang, et al. A survey on evaluation of large language models. *arXiv preprint arXiv:2307.03109*, 2023.
- [6] F. Christianos, G. Papoudakis, and S. V. Albrecht. Pareto actor-critic for equilibrium selection in multi-agent reinforcement learning, 2023. URL <https://arxiv.org/abs/2209.14344>.
- [7] C. S. de Witt, T. Gupta, D. Makoviichuk, V. Makoviychuk, P. H. S. Torr, M. Sun, and S. Whiteson. Is Independent Learning All You Need in the StarCraft Multi-Agent Challenge?, Nov. 2020.
- [8] European Commission. Artificial Intelligence Act. <https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=celex:52021PC0206>, 2021. Accessed: 2024-01-22.
- [9] I. Gabriel. Artificial intelligence, values, and alignment. *Minds and Machines*, 30:411–437, 09 2020. doi: 10.1007/s11023-020-09539-2.
- [10] Z. Gábor, Z. Kalmár, and C. Szepesvári. Multi-criteria reinforcement learning. In *ICML*, volume 98, pages 197–205, 1998. URL <https://sites.ualberta.ca/~szepesva/papers/multi98.ps.pdf>.
- [11] E. Hughes, J. Z. Leibo, M. Phillips, K. Tuyls, E. Dueñez-Guzman, A. García Castañeda, I. Dunning, T. Zhu, K. McKee, R. Koster, et al. Inequity aversion improves cooperation in intertemporal social dilemmas. *Advances in neural information processing systems*, 31, 2018.
- [12] E. Hughes, J. Z. Leibo, M. Phillips, K. Tuyls, E. Dueñez-Guzman, A. García Castañeda, I. Dunning, T. Zhu, K. McKee, and R. Koster. Inequity aversion improves cooperation in intertemporal social dilemmas. *Advances in neural information processing systems*, 31, 2018.
- [13] J. Jumper, R. Evans, A. Pritzel, T. Green, M. Figurnov, O. Ronneberger, K. Tunyasuvunakool, R. Bates, A. Židek, A. Potapenko, et al. Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873): 583–589, 2021.
- [14] M. Lapan. *Deep Reinforcement Learning Hands-On: Apply modern RL methods, with deep Q-networks, value iteration, policy gradients, TRPO, AlphaGo Zero and more*. Packt Publishing Ltd, June 2018.
- [15] J. Z. Leibo, V. F. Zambaldi, M. Lanctot, J. Marecki, and T. Graepel. Multi-agent reinforcement learning in sequential social dilemmas. *CoRR*, abs/1702.03037, 2017. URL <http://arxiv.org/abs/1702.03037>.
- [16] A. Lin. Binary search algorithm. *WikiJournal of Science*, 2(1):1–13, 2019.
- [17] R. Noothigattu, D. Bouneffouf, N. Mattei, R. Chandra, P. Madan, R. Kush, M. Campbell, M. Singh, and F. Rossi. Teaching ai agents ethical values using reinforcement learning and policy orchestration. *IBM Journal of Research and Development*, PP:6377–6381, 09 2019. doi: 10.1147/JRD.2019.2940428.
- [18] C. H. Papadimitriou and T. Roughgarden. Computing equilibria in multi-player games. In *SODA*, volume 5, pages 82–91, 2005.
- [19] G. Papoudakis, F. Christianos, L. Schäfer, and S. V. Albrecht. Benchmarking multi-agent deep reinforcement learning algorithms in cooperative tasks. In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks (NeurIPS)*, 2021. URL <http://arxiv.org/abs/2006.07869>.
- [20] G. Papoudakis, F. Christianos, L. Schäfer, and S. V. Albrecht. Benchmarking Multi-Agent Deep Reinforcement Learning Algorithms in Cooperative Tasks, Nov. 2021. URL <http://arxiv.org/abs/2006.07869>.
- [21] M. O. Riedl and B. Harrison. Using stories to teach human values to artificial agents. In *AAAI Workshop: AI, Ethics, and Society*, 2016.
- [22] M. Rodríguez-Soto, M. López-Sánchez, and J. A. Rodríguez-Aguilar. A structural solution to sequential moral dilemmas. In *Proceedings of the 19th international conference on autonomous agents and multiagent systems*, pages 1152–1160, 2020.
- [23] M. Rodríguez-Soto, M. López-Sánchez, and J. A. Rodríguez-Aguilar. Multi-objective reinforcement learning for designing ethical environments. In *IJCAI*, pages 545–551, 2021.
- [24] M. Rodríguez-Soto, M. López-Sánchez, and J. A. Rodríguez-Aguilar. Multi-objective reinforcement learning for designing ethical multi-agent environments. *Neural Computing and Applications*, pages 1–26, 2023.
- [25] D. M. Roijers. Multi-objective decision-theoretic planning. *AI Matters*, 2(4):11–12, 2016.
- [26] D. M. Roijers, P. Vamplew, S. Whiteson, and R. Dazeley. A survey of multi-objective sequential decision-making. *Journal of Artificial Intelligence Research*, 48:67–113, 2013.
- [27] F. Rossi and N. Mattei. Building ethically bounded ai. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33:9785–9789, 07 2019. doi: 10.1609/aaai.v33i01.33019785.
- [28] S. Russell, D. Dewey, and M. Tegmark. Research priorities for robust and beneficial artificial intelligence. *Ai Magazine*, 36(4):105–114, 2015.
- [29] R. Rădulescu, P. Mannion, D. M. Roijers, and A. Nowé. Multi-objective multi-agent decision making: a utility-based analysis and survey. *Autonomous Agents and Multi-Agent Systems*, 34(1), Apr. 2020. ISSN 1387-2532, 1573-7454. doi: 10.1007/s10458-019-09433-x. URL <http://link.springer.com/10.1007/s10458-019-09433-x>.
- [30] M. Serramia, M. Rodríguez-Soto, M. López-Sánchez, J. A. Rodríguez-Aguilar, F. Bistaffa, P. Boddington, M. Wooldridge, and C. Ansotegui. Encoding ethics to compute value-aligned norms. *Minds and Machines*, 33(4):761–790, 2023.
- [31] J. Skalse, L. Hammond, C. Griffin, and A. Abate. Lexicographic Multi-Objective Reinforcement Learning. In *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence*, pages 3430–3436, Vienna, Austria, July 2022. International Joint Conferences on Artificial Intelligence Organization. ISBN 978-1-956792-00-3. doi: 10.24963/ijcai.2022/476. URL <https://www.ijcai.org/proceedings/2022/476>.
- [32] N. Soares and B. Fallenstein. Aligning superintelligence with human interests: A technical research agenda. *Machine Intelligence Research Institute (MIRI) technical report*, 8, 2014.
- [33] R. Sutton. The bitter lesson. *Incomplete Ideas (blog)*, 13(1):38, 2019.
- [34] E. Tennant, S. Hailles, M. Musolesi, et al. Modeling moral choices in social dilemmas with multi-agent reinforcement learning. In *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence, IJCAI-23*, pages 317–325, 2023.
- [35] A. Tercan and V. S. Prabhu. Thresholded Lexicographic Ordered Multiobjective Reinforcement Learning, Sept. 2024. URL <http://arxiv.org/abs/2408.13493>. arXiv:2408.13493 [cs].
- [36] P. Vamplew, R. Dazeley, C. Foale, S. Firmin, and J. Mummery. Human-aligned artificial intelligence is a multiobjective problem. *Ethics and Information Technology*, 20, 03 2018. doi: 10.1007/s10676-017-9440-6.
- [37] P. Vamplew, C. Foale, R. Dazeley, and A. Bignold. Potential-based multiobjective reinforcement learning approaches to low-impact agents for ai safety. *Engineering Applications of Artificial Intelligence*, 100, 04 2021. doi: 10.1016/j.engappai.2021.104186.
- [38] I. Van de Poel and L. Royakkers. *Ethics, technology, and engineering: An introduction*. John Wiley & Sons, 2023. URL <https://books.google.com/books?hl=en&lr=&id=SYq4EAAAQBAJ&oi=fnd&pg=PR10&dq=Ethics,+technology,+and+engineering:+an+introduction.&ots=O46iJck9zx&sig=ZDX7DwrsWkNNBWMe0YISgqlv2hY>.
- [39] O. Vinyals, I. Babuschkin, W. M. Czarnecki, M. Mathieu, A. Dudzik, J. Chung, D. H. Choi, R. Powell, T. Ewalds, P. Georgiev, et al. Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature*, 575(7782):350–354, 2019.
- [40] B. Von Stengel. Computing equilibria for two-person games. *Handbook of game theory with economic applications*, 3:1723–1759, 2002.
- [41] Y.-H. Wu and S.-D. Lin. A low-cost ethics shaping approach for designing reinforcement learning agents. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- [42] P. R. Wurman, S. Barrett, K. Kawamoto, J. MacGlashan, K. Subramanian, T. J. Walsh, R. Capobianco, A. Devlic, F. Eckert, F. Fuchs, et al. Outracing champion gran turismo drivers with deep reinforcement learning. *Nature*, 602(7896):223–228, 2022.
- [43] C. Yu, A. Velu, E. Vinitzky, J. Gao, Y. Wang, A. Bayen, and Y. Wu. The surprising effectiveness of ppo in cooperative multi-agent games. *Advances in Neural Information Processing Systems*, 35:24611–24624, 2022.
- [44] H. Yu, Z. Shen, C. Miao, C. Leung, V. R. Lesser, and Q. Yang. Building ethics into artificial intelligence. In *IJCAI*, page 5527–5533, 2018.