ALTARED ENVIRONMENTS: THE ROLE OF NORMA-TIVE INFRASTRUCTURE IN AI ALIGNMENT

Anonymous authors

Paper under double-blind review

ABSTRACT

Cooperation is central to human societies, which they achieve by constantly tackling the alignment problem of ensuring self-interested individuals act in ways that benefit the groups in which they live. As AI agents become pervasive in shared environments, it will be similarly crucial for them to align with the cooperative goals of human groups. Current AI alignment research largely focuses on embedding specified or learned norms into agents to achieve this cooperation. While valuable, this approach overlooks the role that institutions play in aligning human behavior to achieve cooperative gains and thus overlooks a potential alignment technique for AI agents. We address this gap by proposing Altared Games, a novel formal extension of Markov games that incorporates an *altar*—a classification institution providing explicit normative guidance to agents. Our approach focuses on a challenging setting where norms are dynamic, thereby requiring agents to adapt to the evolving norm content represented by the altar. Using multi-agent reinforcement learning (MARL) as a computational model of AI agents, we conduct experiments in two mixed-motive environments: Commons Harvest, which models resource sustainability, and Allelopathic Harvest, which involves coordination under conflicting incentives. Our results demonstrate that the altar enables agents to adapt effectively to dynamic norms, engage in accurate sanctioning, and achieve higher social welfare compared to systems without a classification institution. These findings highlight the importance of normative institutions in fostering cooperative, adaptable AI agents operating in complex real-world settings.

031 032

033

004

010 011

012

013

014

015

016

017

018

019

021

023

024

025

026

027

028

029

1 INTRODUCTION

The alignment challenge – how to ensure that self-interested individuals adopt behaviors that benefit the groups they live in – has been a persistent and evolving problem throughout human history. Efforts to address this challenge have been foundational for achieving the kind of ultra-cooperative societies humans have built. By engaging in schemes of task specialization, exchange, and mutual aid, which require individuals to follow group norms of appropriate behavior, humans have achieved levels of cooperation far beyond anything we see in other mammals Henrich (2016).

040 Integrating AI into human society extends the alignment challenge to artificial agents: how do we 041 make sure these agents take actions that align with the norms of appropriate behavior that undergird 042 our complex cooperative schemes? To date, this challenge has largely been framed in terms of how 043 we embed values and norms into AI systems. And while this approach has been important for the 044 safe deployment of existing systems, it is inherently limited. Human norms are dense-just about everything we say or do is subject to normative evaluation. Norms are dynamic, constantly adapting to changes in environments, populations, and information, and open to continuing contestation. 046 And norms are highly differentiated: they range from ineffable standards such as how long it is 047 appropriate to make eye contact with a stranger to legible norms about color code to at a funeral or 048 how much food to take from a shared plate to formal legal requirements such as the obligation to take reasonable care while driving or to put away your garbage cans within 24 hours of collection. It is simply not possible to articulate all our values and norms Hadfield-Menell & Hadfield (2018). 051

Tackling the AI alignment challenge in a robust way will require taking fully on board the density,
 dynamism and differentiation of human norms and the lessons from how human groups throughout
 history have tackled the alignment challenge, successfully enough to have achieved extraordinary



064

Figure 1: Overview of the **Altared Games**: Markov Game extension including third party enforcement mechanism and an environmental feature, called *altar*, encoding the norm that evolve over time. Learning unfolds in this order: Agents first learn to punish in accordance with the hidden reward structure for norms. In the presence of altar, they learn to map the observation to this hidden structure. This enables them to potentially predict the sanctioning behavior of other agents (a challenging problem) and then, as a consequence, they learn to avoid sanctions and comply with norms.

gains from cooperation. What we see in the human model is that humans do not, by and large, 071 encode specific values and norms early in life that then guide lifelong behavior. Instead, human 072 societies rely on constantly evolving interactions with other agents, normative signals, and, most 073 fundamentally, normative institutions - authoritative common knowledge structures like groups of 074 elders or courts that articulate, interpret and adapt norms-to align individual behaviors in dynamic 075 environments and with dynamic populations. We can understand the emergence and evolution of 076 normative institutions as a response to the density and dynamism of norms. Individuals in groups 077 that lack an authoritative normative institution must extract information about the current norms and how they are being interpreted and enforced by other agents from agent behavior alone. Individuals in groups with an authoritative normative institution face a less computationally taxing and error-079 laden challenge in maintaining coordination of their enforcement and compliance behavior with the group. This generates group benefits in the form of increased social stability. 081

082 In this paper, we draw on this human model to investigate whether a normative institution can im-083 prove the capacity of architecturally simple AI agents to adapt to dynamic norms while solving social dilemmas Dawes & Messick (2000). We begin with the theoretical framework of (Hadfield & 084 Weingast, 2012), which introduced a rational agent model of normative social order. In this frame-085 work, the challenge of inducing behaviors aligned with a group's norms resolves to the challenge of incentivizing and coordinating agents to punish in accordance with the normative classification 087 (which behaviors are allowed, which are not allowed) articulated by a public classification insti-880 tution. This framework provides a microfoundational account of the decentralized enforcement 089 mechanisms seen in human societies, such as social disapproval or exclusion, that are a primary mechanism for incentivizing norm compliance. 091

We adapt this theoretical framework to the AI context with an extension to Markov games – a robust computational method for modeling sequential decision-making in multi-agent environments. We formalize decentralized enforcement in this setting by endowing agents in a multi-agent reinforcement learning (MARL) setting (following Perolat et al. (2017)) with a sanctioning technology by which they can deliver costly punishment to other agents. Following Köster et al. (2022), we implement norms by encoding rewards for agents that use their sanctioning technology to punish agents that have taken actions deemed by the norm to be punishable. If all or most agents reliably punish in this way, they also avoid violating the norm themselves to avoid punishment from others.

099 The punishment rewards in this framework are "hidden" in the sense that they are supplied by the 100 environment and not modeled by the agents. How to earn rewards for punishing is thus a learning 101 problem, one that (Köster et al., 2022) show MARL agents can solve with static norms/rewards. We 102 make this learning problem harder and more realistic by implementing dynamic norms that change in 103 a randomly controlled way during training. We then evaluate the impact of introducing a normative 104 classification institution, a feature that provides a publicly observable representation of the current 105 norm, that is, the current reward structure for punishment. We call this feature an **altar**, to capture the idea of an authoritative focal point in the environment that articulates a group's shared norms 106 or laws, and propose a novel extension of the Markov game setting called Altared Games. Our 107 research question then is: does the introduction of an altar make it easier for agents to adapt their

punishment behaviors to changes in the norms and thus for a group of agents to maintain dynamic normative social order?

To empirically investigate this framework, we conduct experiments in two mixed-motive games 111 with different cooperative challenges: Commons Harvest, modeling resource sustainability under 112 the tragedy of the commons; and Allelopathic Harvest, involving equilibrium selection and a free-113 riding challenge. To address our research question we use a methodology of controlled hypothesis 114 testing to isolate the role of the altar, ensuring that the effects of the altar can be assessed indepen-115 dently by holding all other factors constant across experimental setups. Specifically, we train agents 116 under three experimental conditions: (1) a vanilla baseline in which agents possess a sanctioning 117 technology but sanctioning only generates private rewards (such as removing a competitor from a 118 contested resource) meaning there are no pro-social rewards for sanctioning and hence no norms; (2) a hidden-rule environment in which agents are rewarded by the environment for punishing in 119 accordance with the current norm; and (3) the altared environment which enriches the hidden-rule 120 environment with an altar that represents the state of the current norm. In both the hidden-rule and 121 the altared conditions the norms follow the same dynamic evolution and in both cases the group 122 of agents would achieve the highest possible payoff if agents immediately shifted their behavior to 123 align with the current norm. (Note that this means we are not in this paper investigating how a group 124 identifies and adopts the optimal norms). Thus a comparison of performance between agents trained 125 with a hidden-rule and with an altar isolates the impact of the altar alone. Our results show that the 126 altar significantly improves agents' ability to adapt to dynamic norms, engage in correct sanctioning 127 behaviors, and achieve higher social welfare efficiently, even under uncertainty. 128

In sum, our work provides a first step toward understanding how normative institutions can enhance alignment in multi-agent systems. Focusing on dynamic norms and explicit institutional guidance, we aim to pave the way for future research into scalable, adaptable, and socially aligned AI systems.

131 132

129

130

2 PRELIMINARIES

133 134

This section establishes the theoretical and formal foundations of this work. We summarize the theoretical framework of normativity, focusing on the role of classification institutions and enforcement mechanisms in sustaining cooperation. We then describe Markov games and its extension to sanction-augmented Markov games, which incorporate third-party enforcement into a computationally rich multi-agent setting. An extended related work discussion is available in Appendix B

140 141

2.1 THEORETICAL FRAMEWORK: HADFIELD-WEINGAST MODEL

142 Our investigation is grounded in a parsimonious rational agent model of normative social order in-143 troduced by Hadfield and Weingast (Hadfield & Weingast, 2012). This model provides a structured 144 perspective on how groups sustain cooperation by leveraging two essential components: a classifi-145 cation institution and an enforcement mechanism. The classification institution provides common 146 knowledge binary classifications of behaviors as either "punishable" or "not punishable," potentially 147 through the application of general principles to specific cases. These classifications reduce ambi-148 guity, creating a shared understanding of acceptable behavior within the group. The enforcement mechanism incentivizes agents to align with these classifications by imposing penalties on punish-149 able actions, encouraging agents to favor "not punishable" behaviors. A stable normative social 150 order is achieved when most agents are mostly in compliance and avoiding punishment. 151

(Hadfield & Weingast, 2012) focus in particular on the case, which describes most of human history and much of modern life as well, in which punishment is decentralized, that is, primarily delivered by ordinary agents (rather than specialized enforcers Hadfield & Weingast (2013).) Agents must therefore be incentivized and coordinated to engage in costly third-party punishment (which could be relatively mild, such as criticism, or more harsh, such as exclusion from the group) and to condition such punishment actions on a shared classification institution..

Although shared classification could be entirely emergent and informal¹, groups that converge on a single authoritative (more formal) classification institution–such as a chief, a group of elders,

160 161

¹There is no entity that tells members of the group that they should honk at a car that is failing to facilitate merging on the highway but everyone in the group could reliably say that this is the norm Bicchieri (2005)

or a court-can enjoy significant benefits. These formalized systems help resolve ambiguities and
 improve normative clarity, provide consistency and help maintain cooperation even in the face of
 dynamic environments and populations (Hadfield, 2017).

166 2.2 MARKOV GAMES

168 Markov games, also known as stochastic games, extend Markov decision processes to multi-169 agent settings, providing a general framework for modeling dynamic interactions among agents. 170 A Markov game is defined as a tuple: (S, A, P, R, γ, n) , where, S is the shared state space, representing all possible configurations of the environment; $A = A_1 \times A_2 \times \cdots \times A_n$ is the joint action 171 space, where A_i denotes the set of actions available to agent i, and n is the total number of agents; 172 $P(s' \mid s, a)$ is the transition function, specifying the probability of transitioning to state s' from state 173 s given the joint action a; $R = (R_1, R_2, \ldots, R_n)$ represents the reward functions for each agent, 174 where $R_i(s, a)$ determines the reward received by agent i after taking action a in state s; $\gamma \in [0, 1)$ 175 is the discount factor, controlling the relative importance of future rewards. 176

In **partially observable** settings, each agent i does not have access to the full state s but in-177 stead receives an observation $o_i \in O_i$, drawn from the observation function $O_i(s)$. The Markov 178 game is extended to a partially observable Markov game (POMG) by redefining the tuple as: 179 $\langle S, A, P, R, \gamma, n, O \rangle$, where $O = (O_1, O_2, \dots, O_n)$ defines the observation spaces of the agents. 180 The addition of partial observability introduces complexity, as agents must infer hidden state in-181 formation from their observations to make optimal decisions. At each timestep, agents observe s 182 (or o_i in the partially observable case), select actions $a_i \in A_i$, and transition to a new state s' 183 based on $P(s' \mid s, a)$. Each agent's goal is to learn a policy $\pi_i : S \to A_i$ (or $\pi_i : O_i \to A_i$) 184 in the partially observable case) that maximizes its expected cumulative discounted reward: $\pi_i^* =$ 185 $\arg \max_{\pi_i} \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t R_i(s_t, a_t)\right]$. A subset of Markov games, mixed-motive settings, involves a combination of cooperative and competitive incentives. These settings model scenarios where agents must balance individual objectives with the collective good, often facing dilemmas such as 187 the equilibrium selection problem, where multiple equilibria, generally with different individual and 188 aggregate payoffs, exist; the free-rider problem, where agents benefit from shared resources without 189 contributing to their production or maintenance; or the tragedy of the commons, where uncoordi-190 nated actions lead to the depletion of shared resources. Studying such settings is central to our 191 investigation, as they highlight the challenges of aligning individual incentives with group goals and 192 provide a rich domain for exploring the role of norms and enforcement mechanisms. 193

Multi-Agent Reinforcement Learning (MARL) provides the computational framework for solv-194 ing Markov games, where agents interact with the environment and each other to optimize their 195 policies. Formally, each agent i learns a policy π_i to maximize its cumulative discounted re-196 ward: $\pi_i^* = \arg \max_{\pi_i} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R_i(s_t, a_t) \mid \pi_1, \dots, \pi_n \right]$, where $a_t = (a_{1,t}, a_{2,t}, \dots, a_{n,t})$ rep-197 resents the joint action at timestep t. MARL approaches can involve optimizing a joint policy $\pi = (\pi_1, \ldots, \pi_n)$ under shared information or decentralized policies where agents act indepen-199 dently. For this work, we consider a simple, practical and scalable approach to MARL is Indepen-200 dent Proximal Policy Optimization (IPPO) de Witt et al. (2020a), a decentralized method where each 201 agent optimizes its policy independently using a variant of Proximal Policy Optimization (PPO).

202 203

204

2.3 SANCTION-AUGMENTED MARKOV GAMES (SMG)

205 We consider an extension of Markov games to include sanctions (punishment), called Sanction-206 Augmented Markov Games (SMG). This framework formalizes how agents can impose penalties on others to enforce compliance with norms. An SMG is defined as: $\langle S, A', P, R, \gamma, n, \Sigma, C, I \rangle$, 207 where S, A, P, R, γ, n retain their meanings from the standard Markov game definition. Each agent 208 *i* has an extended action space $A'_i = A_i \cup \{\sigma\}$, where σ is a common sanctioning action (e.g. 209 zapping or criticism) available to all agents. The transition dynamics $P(s' \mid s, a)$ determine the next 210 state s', while the reward functions $R = (R_1, R_2, \dots, R_n)$ incorporate the effects of sanctions into 211 individual incentives. This form of SMG (not including C and I, defined below) was introduced in 212 Perolat et al. (2017) where they endowed agents with a punishment (sanctioning) technology. 213

Köster et al. (2022) took this extension a step further by implementing a **hidden classification rule** that implemented norms by rewarding agents for sanctioning behaviors designated exogenously (i.e. by the researchers) as norm violations. Formally, the sanction cost function $C = (C_i, C_j)$ defines 216 the costs associated with sanctions: $C_i(\sigma) = p_i$, the cost incurred by the sanctioning agent i, and 217 $C_i(\sigma) = p_i$, the cost (penalty) incurred by the sanctioned agent j. The indicator variable then 218 implements a classification of actions as either norm violations or not: $I_i(s, a) = 1$ if the action a 219 taken previously by agent j violates a norm, and $I_i(s, a) = 0$ otherwise. To incentivize sanctioning 220 of norm violations, the sanctioning agent receives a reward q that offsets the cost p_i , resulting in a net positive reward $(q - p_i)$ if the sanctioned agent j violated a norm $(I_i(s, a) = 1)$ in the prior step. 221 We call sanctioning of designated norm violations *correct sanctioning*. Conversely, if the sanctioned 222 agent did not violate the norm $(I_i(s, a) = 0)$, the sanctioning agent incurs the full cost of sanctioning 223 p_i without any offsetting reward. We call this *incorrect sanctioning*. Formally, the reward for the 224 sanctioning agent i is: $R_i(s, a, \sigma) = R_i(s, a) - C_i(\sigma) + I_i(s, a) \cdot q$. For the sanctioned agent j, the 225 reward function reflects the penalty for being sanctioned, irrespective of whether the agent violated 226 a norm: $R_i(s, a, \sigma) = R_i(s, a) - C_i(\sigma)$. This formalization captures both the costs and rewards 227 of sanctions, emphasizing the role of accuracy and cost-effectiveness in enforcement. Sanctioning 228 agents are incentivized to sanction correctly to offset their costs, while incorrect sanctions lead to 229 a net penalty. Sanction-augmented environments align well with the (Hadfield & Weingast, 2012) 230 theoretical framework, modeling third-party enforcement to secure normative social order.

231 232

3 OUR APPROACH: ALTARED GAMES

233 234

As discussed in 2.3, (Köster et al., 2022) leveraged the SMG framework for introducing a hidden classification rule, rewarding agents for sanctioning behaviors aligned with predefined but implicit norms. They demonstrated that MARL agents can efficiently learn to punish, and therefore comply, with researcher-set norms. In this paper we extend this framework in two new ways:

First, we address the challenge of **dynamic norms**, where the classification behaviors as punishable or acceptable evolves over time. Second, we introduce a **normative institution**, called the **altar**, which encodes the prevailing norms in the environment. The altar is implemented as an observational feature of the environment and does not modify the structure of the underlying SMG. The altar makes normative content legible to agents. We hypothesize that as a result of this enrichment of the environment, the altar facilitates agent learning and coordination in environments with dynamic norms relative to the hidden rules environment studied by (Köster et al., 2022).

246

247 3.1 DYNAMIC NORMS 248

249 In the context of Sanction-Augmented Markov Games (SMGs), we formalize dynamic norms as 250 a time-dependent mapping: $N_t: S \to \mathcal{A}$, where $N_t(s) \subseteq \mathcal{A}$ defines the set of acceptable (not 251 punishable) actions in state s at time t. The evolution of norms is governed by an update function: $N_{t+1} = f(N_t, \Phi)$, where f captures the mechanism of norm evolution, and Φ represents triggers or 252 drivers of change. We do not model the determinants of norm evolution but these drivers could be 253 thought of as arising from external inputs (e.g., regulatory updates or environmental changes), agent-254 driven mechanisms (e.g., collective decision-making or voting), or stochastic events (e.g., resource 255 depletion or unexpected disturbances). 256

Dynamic norms pose two main challenges for agents. First, dynamic norms require agents to con-257 tinuously track the norm as it evolves, updating their internal models based on observed rewards, 258 sanctions, and environmental cues. Second, agents must adjust their strategies to align with shifting 259 expectations about rewards while navigating a mixed-motive setting. In our setup, these challenges 260 are particularly acute. The immediate impact of a change in the norm is not on the rewards asso-261 ciated with actions that either comply or not with the norm; rather, it is on the rewards associated 262 with sanctioning actions. The impact of norm change on compliance is only derivative: if enough 263 agents adapt their sanctioning behaviors to accord with the new norm, then agents will adapt their 264 compliance behaviors to accord with the new norm. ((Köster et al., 2022) show that this learning 265 process is sequential: MARL agents first learn to punish in accordance with the hidden reward struc-266 ture for norms and then, as a consequence, they learn to comply with norms. Predicting the rewards 267 associated with compliance and non-compliance (which can impact individual payoffs as norms in a mixed-motive setting generally will sometimes require agents to forego self-interested actions in 268 favor of pro-social actions), then is a very difficult problem in a multi-agent setting as it requires 269 predicting the enforcement behavior of other agents.

Real-world examples include resource management scenarios, where norms shift in response to scarcity, and traffic systems, where acceptable behaviors adapt to changing infrastructure or population density. By incorporating dynamic norms into SMG, we aim to model these complexities and investigate how agents operate in environments with evolving expectations.

274

275 276 3.2 Altared Games

277 Our aim is to test the value of a normative institution, which we call the altar and which encodes 278 the prevailing reward structure for punishment (the norms). We thus further extend the SMGs 279 framework to incorporate the altar feature and call this Altared Sanction-Augmented Markov Games 280 (Altared SMGs), which we will refer to as Altared Games for short. An Altared SMG is defined as: 281 $\langle S, A', P, R, \gamma, n, \Sigma, C, I, \mathcal{M}_{altar} \rangle$, where $S, A', P, R, \gamma, n, \Sigma, C$, and I retain their meanings from 282 the SMG framework, and $\mathcal{M}_{altar}: S \times \mathcal{A} \to \{0,1\}$ is a mapping function managed by the envi-283 ronment, encoding the normative classification of actions. It specifies whether an action a in state 284 s complies with the norm, with $\mathcal{M}_{\text{altar}}(s, a) = 1$ indicating compliance and $\mathcal{M}_{\text{altar}}(s, a) = 0$ indi-285 cating violation. Combined with the sanctioning mechanism described above, the altar thus encodes the sanctioning reward structure, indicating what constitutes correct sanctioning. 286

287 Agents do not have direct access to \mathcal{M}_{altar} . Instead, when visiting a designated subset of states 288 $S_{\text{altar}} \subseteq S$, they receive an observation o_{altar} that implicitly reflects the normative content encoded 289 by $\mathcal{M}_{\text{altar}}$. The indicator variable $I_j(s, a)$, which specifies whether the action a by agent j violates 290 the norm, is derived implicitly from \mathcal{M}_{altar} : $I_j(s, a) = 1 - \mathcal{M}_{altar}(s, a)$. Thus, while \mathcal{M}_{altar} governs the normative structure of the environment, agents must infer this structure through observation and 291 feedback. Moreover, the reward functions for the sanctioning agent i, $R_i(s, a, \sigma)$ and the sanctioned 292 agent j, $R_i(s, a, \sigma)$ remain consistent with the SMG framework. The altar observations merely 293 provide agents with additional information about the rewards for sanctioning. Thus, while there is no direct cost associated with visiting S_{altar} , interactions with the altar involve implicit opportunity 295 costs: Agents forgo potential reward-generating actions during the time spent visiting S_{altar} . 296

This formalization bridges the gap between implicit norm enforcement in hidden rule systems and explicit norm representation. By incorporating the altar into the SMG framework, we create a testbed for investigating how observable classification institutions influence agent learning, coordination, and compliance in dynamic, multi-agent environments.

301 302

4 EXPERIMENTS

303 304

The objective of our experiments is to evaluate the impact of the altar, on agent behavior, norm learning, enforcement, and compliance in dynamic multi-agent systems. Specifically, we aim to understand whether making norms observable through the altar improves agents' ability to align with evolving norms, enforce compliance, and achieve higher overall system efficiency compared to configurations without explicit institutional representation.

To explore these questions, we use two mixed-motive environments: Commons Harvest Pero-310 lat et al. (2017) and Allelopathic Harvest Köster et al. (2020). To realize these environments, we 311 leverage the Melting Pot Suite Agapiou et al. (2023); Leibo et al. (2021), a flexible research platform 312 that provides high-fidelity multi-agent environments with diverse incentive structures and interde-313 pendencies. Its extensibility allows us to adapt these environments systematically to include explicit 314 institutional mechanisms like the altar. The aim of our experiment is to test the hypothesis that 315 the presence of an altar that encodes norms improves the capacity of agents to implement norms. 316 For this reason, following (Köster et al., 2022), we exogenously control the content of norms. We 317 test our hypothesis by training agents under three experimental conditions. In the Vanilla Base-318 line (Markov Game), the original Markov game is used without norms or sanctioning technology, 319 and agents maximize individual rewards without external guidance or sanctions. This gives us a 320 reference point to assess the group benefits achieved if the agents are able to implement our deliber-321 ately group-beneficial norms. Our second condition, Hidden Rule SMG, introduces the sanctioning technology and the hidden reward structure that rewards sanctioning according to the current norms. 322 Finally, the **Altared SMG** condition incorporates the altar but is otherwise the same as the Hidden 323 Rule SMG condition, with dynamic norms that follow the same evolution and the same rewards for punishing according to these norms. By comparing performance for agents trained under these three conditions, we aim to isolate the effects of the altar on agent behavior and system outcomes.

4.1 Environments: Core, SMG and Altared Versions

In this section, we provide an overview of the two environments used in our experiments: **Commons Harvest** and **Allelopathic Harvest**. For each environment, we first outline the core mechanics, describing the resource dynamics and agent interactions that define the setting. We then explain how the Sanction-Augmented Markov Game (SMG) version of the environment is constructed, leading to the hidden rule mechanism for enforcing norms. Finally, we detail the steps taken to convert these environments into their **Altared** versions, explicitly incorporating the altar as an observable institution encoding norms.

336 337

327

328

4.1.1 COMMONS HARVEST².

Core Mechanics. In this environment, agents aim to collect apples scattered across six distinct patches, earning a reward of +1 for each apple consumed. Apple regrowth depends on the density of neighboring apples within a Euclidean radius of 2, with probabilities decreasing as local density declines: 0.025 for three or more neighbors, 0.005 for two, 0.001 for one, and 0 for none. Overharvesting a patch depletes it permanently, requiring agents to reduce collection to sustain resources. A social dilemma ensues: consuming the last apple in a patch generates individual rewards but risks permanent patch depletion, leading to the tragedy of the commons.

For this work, we divide the six patches into three zones: the top two patches are designated red, the middle two are blue, and the bottom two are green. Agents are initially gray, but take on the color of the zone from which they consume apples. This color change makes their collection behavior observable by other agents. This setup lays the groundwork for introducing sanctioning based on zones in subsequent versions. We then train agents in three conditions described in Appendix A.1.

Achieving normative alignment in this environment translates to agents adapting to evolving norms and sanctioning other agents correctly. Because the norm is adjusted to reflect the current supply of apples across different patches, correct sanctioning in accordance with current norms incentives agents to adapt their harvesting behavior to the health of apple supply, mitigate overharvesting and thereby achieving higher collective welfare over time.

355 356

4.1.2 Allelopathic Harvest³

357 **Core Mechanics:** This environment poses both the coordination and the free-rider problem, making 358 it challenging for agents to reach a welfare maximizing outcome. Specifically, in this environment, 359 there are berries of three different colors and sixteen agents can plant and consume berries. Agents 360 get reward for consuming any colored berry (+1) but receive higher reward for consuming their 361 preferred color berry (+2). Planting does not generate any reward or cost and hence agents have 362 no direct incentive to plant., leading to a free-rider problem. The agents can only consume ripened 363 berries and the berry ripening rate is directly proportional to the fraction of the largest amount of berry color. Hence, if all three colors are equally distributed, berries will have the slowest ripening 364 rate and achieving a monoculture of a single berry color will generate the highest berry ripening 365 rate, thereby giving a chance to agents to accumulate more reward (equilibrium selection problem). 366 Agents are initially gray, but take on the color of the berry they plant. This color change makes 367 their collection behavior observable by other agents. This setup lays the groundwork for introducing 368 sanctioning based on berry color for which monoculture is desired. We then train agents in three 369 conditions discussed in Appendix A.2. 370

Achieving normative alignment in this environment translates to agents adapting to evolving norms and sanctioning other agents correctly. Because the norm is adjusted to reflect the currently desired

³⁷³

 ² Perolat et al. (2017) introduced this environment to investigate the ability of multi-agent reinforcement learning agents to coordinate in solving common-pool resource appropriation problems, building on the mechanics first outlined in Janssen et al. (2010)

 ³ Köster et al. (2020) introduced this environment to investigate the ability of multi-agent reinforcement
 learning agents to overcome free-rider problem while solving equilibrium selection problem rooted in the al lelopathic mechanic, previously studied in in Leibo et al. (2019)



Figure 2: Results on Altared Commons Harvest: Adjusted reward mean discounts the reward obtained for sanctioning. All experiments run for 5 seeds.

monoculture color, correct sanctioning in accordance with current norms incentives agents to adapt their planting behavior to the desired monoculture while avoiding to free-ride, thereby achieving higher collective welfare over time.

4.2 Results

386

387 388 389

390

391

392 393

394

395 Our empirical investigation focuses on assessing the impact of the altar on agents' capacity to im-396 plement norms, compared to environments without the altar. During training, we expect the learning 397 process to unfold as follows: agents first learn to recognize nonacceptable behaviors by receiving 398 rewards for sanctioning violations, enabling them to enforce punishments correctly. Over time, this 399 enforcement leads agents to predict actions likely to result in sanctions, prompting them to learn 400 compliant behavior by avoiding such actions. This progression drives agents toward maintaining a normative social order. In the presence of the altar, agents would be required to visit it periodically 401 to update their understanding of the prevailing norm. They must learn to map altar observations 402 to appropriate sanctioning behaviors, potentially facilitating faster and more accurate adaptation to 403 changing norms. We present our results through the lens of this learning process, comparing agent 404 performance in environments with and without the altar (everything else being the same) at each 405 stage of this progression. The results are reported over 5 seeds for all experiments and further train-406 ing details are available in Appendix C 407

We highlight that agents not engaged in a normative system face prohibitive difficulty in learning the restraint required in Commons Harvest to avoid the tragedy of the commons, as observed in agents trained under the vanilla condition. In the Allelopathic Harvest environment, these agents tend to free-ride by consuming berries indiscriminately, preventing any increase in the growth rate of berries and resulting in stagnation at a specific reward level.

Agents learn correct sanctioning behavior in the presence of an altar. As a first result, we mea-413 sure the impact of institution on the ability of agents to enforce punishments correctly. For this, 414 we plot the fraction of the correct and incorrect sanctions that agents engage in over the course of 415 training for the baseline without institution and our approach. Figures 2b, 2c and Figures 3b, 3c 416 shows the results for Common Harvest and Allelopathic Harvest environments respectively. In the 417 Commons Harvest environment, agents trained in the Altared SMG framework quickly learn to per-418 form the majority of their sanctions correctly and adapt to dynamic norms. They maintain a high 419 fraction of correct sanctioning behavior over extended training periods while significantly reducing 420 variance compared to the Hidden Rule SMG baseline. This highlights the advantage provided by 421 the altar feature in facilitating agents' understanding of the sanctioning reward structure. In the Al-422 lelopathic Harvest environment, agents trained in Altared SMG initially struggle to identify correct sanctioning behaviors. However, they eventually match the performance of the Hidden Rule SMG 423 baseline and, over time, appear to surpass it. Additionally, the variance in performance is consis-424 tently lower for Altared SMG compared to Hidden Rule SMG, underscoring the stabilizing effect of 425 explicit institutional guidance. 426

Agents learn to visit the altar consistently. In the challenging setup of dynamic norms, it is crucial
 for agents to learn to visit the altar at regular intervals to stay updated on the evolving normative
 content. To evaluate this behavior, we tracked the number of visits made by agents to the altar over
 the course of training. Figures 4b and 4a depict the visitation patterns for the Commons Harvest
 and Allelopathic Harvest environments, respectively. Our results show that, over time, agents in both
 environments converge on a consistent visitation pattern, maintaining a stable frequency of visits to



(a) Altar visits in Allelopathic Harvest

vest

(c) Cumulative Depletion rate

Figure 4: Results on Altar visits in both environments and Depletion rate in Commons Harvest

the altar. This behavior suggests that agents effectively learn the importance of periodic updates 455 from the altar to adapt to the dynamic norms. 456

457 Agents obtain high collective welfare, more efficiently when trained in the presence of altar. 458 It is important to note that agents in both the Altared SMG and Hidden Rule SMG setups receive 459 rewards for correctly sanctioning violations. However, this reward is artificial and is solely intended 460 to train agents to learn proper sanctioning behavior. Sanctioning is inherently costly for both the source and target agents, and it is only justified if the rewards obtained from the base environment 461 outweigh the associated costs. To account for this, we report the Adjusted Mean Reward in our 462 results, which excludes the rewards earned from correct sanctioning. This metric ensures a fair 463 evaluation of overall performance by focusing on the net benefits derived from the base environment 464 while still incorporating the costs associated with sanctioning. 465

As shown in Figure 2a, Altared SMG demonstrates strong per-466 formance in Commons Harvest environment, where agents attain 467 higher reward quickly and are able to sustain the increase in their 468 collective reward. This requires the agents to show collective re-469 straint towards harvesting apples from the zones that are nearing 470 depletion. Surprisingly, the agents in the strong Hidden Rule SMG 471 are not able to learn to show restraint and are not able to attain 472 higher reward. It is important to note that the only difference be-473 tween the baseline and our work is the presence of the institution in 474 the environment. The rationale behind this performance is the crux 475 of our position that institutions are important tools to achieve col-476 lective alignment – institutions take away the burden of enormous 477 amount of computation required by the agents in order to understand the normative social order and reason about it. This effect 478 results in the overall reduction in coordination costs, thereby im-479



Figure 5: Visits to the altars representing content not correlated with sanctioning rewards

proving the efficiency of achieving cooperative outcomes. Qualitatively, one of the most important 480 measure of success in Commons Harvest environment is the ability of the agents to foster sustain-481 ability. To assess this, we compute the cumulative depletion rate of apples that the agents cause 482 over the training time. Figure 4c showcases the ability of agents trained in presence of altar to 483 significantly reduce the depletion rate. 484

In the Allelopathic Harvest, Figure 3a similarly good results, albeit only marginally better than 485 the Hidden Rule SMG baseline. We note that while the difference looks small in the plot, the

453 454

486 Altared SMG is better by a score of 200 points and with much less less variance compared to 487 the Hidden Rule SMG (> 200 vs 60). This demonstrates that the institution already provides the 488 necessary information to learn to attain and maintain higher welfare quicker and more reliably. In the 489 Appendix D, we discuss a non-dynamic version of this environment to perform several qualitative 490 analysis tasks which gives more insights into the behavior of agents in the presence of altar in this environment. Altared SMG's marginal performance improvement over Hidden Rule SMG in the 491 current results can be attributed to the artifact that the training runs on Allelopathic harvest have not 492 finished and it appears that agents have just begun to learn correct sanctioning behavior which may 493 be the limiting factor. <u>191</u>

Altar provides value in terms of system robustness as agents learn to respond to different institutional configurations

To address this question, we consider two variations over the basic version of Altared Allelopathic
 Harvest environment focusing on the robustness of the system when faced with different institutional
 configurations. More details on these environments are available in Appendix E.1.2.

Alleopathic Harvest Limited. This has same altar dynamics as Alleopathic Harvest Altar except that the institution is visible to the agents for only few time steps after the color change mimicking how some institutions are only accessible at particular times. In our results, we observe that the agents start visiting the institution during the first half of the interval since the color changes, thereby continuing to continue learning to adapt the dynamics institutional content despite limited visibility.

Allelopathic Harvest Conflict. This has two extra altars in the environment which serves as distractors to the agent. These altar will display information that do not align with the central altar. But the central altar is the one with the correct prescription and hence the agents need to learn to decrease their visit to the incorrect altar to be able to keep improving their ability to cooperate. Figure 5 demonstrates that the agents indeed learn to cut down their visits to the incorrect altar significantly and keep improving their overall reward, thereby responding and adapting to the correct institution.

511 512

5 DISCUSSIONS AND CONCLUDING REMARKS

513

514 In this work, we draw on the model of human societies to investigate whether a normative institu-515 tion can improve the capacity of architecturally simple AI agents to adapt to dynamic norms while 516 solving social dilemmas. Building on the theory of rational agent model, we propose a formal ex-517 tension of Markov games, called Altared Games, which focuses on the decentralized enforcement 518 mechanisms in multi-agent systems and introduces a feature called an **altar**, hat provides a publicly 519 observable representation of the current norm, that is, the current reward structure for punishment. 520 Using multi-agent reinforcement learning, we examine whether the introduction of an altar make it easier for agents to adapt their punishment behaviors to changes in the norms and thus for a group 521 of agents to maintain dynamic normative social order. In a modified Allelopathic Harvest game 522 and Commons Harvest game, we perform a controlled hypothesis driven testing and demonstrate 523 superior performance of agents trained with an altar compared to those without it. While Altared 524 Games have a strong theoretical grounding, it is also highly intuitive - institutions reduce the cogni-525 tive burden on agents when addressing cooperation challenges by tracking key elements that sustain 526 and promote normative social order. This enables agents to efficiently engage in coordinated, co-527 operative behaviors, facilitating quicker and more effective collective action. This work provides 528 a first step toward understanding how normative institutions can enhance alignment in multi-agent 529 systems. By focusing on dynamic norms and incorporating explicit institutional guidance, we aim 530 to pave the way for future research into scalable, adaptable, and socially aligned AI systems

531 Our current exposition provides a recipe for designing environments and systems with normative 532 infrastructure as a key component, we strongly believe that this direction is ripe with immediate 533 future avenues. We posit that our approach will be particularly effective in promoting generaliza-534 tion, adaptability, and robustness across environments with different normative institutions. Agents that learn to recognize altar information and correlate enforcement behaviors with norm content will 536 adapt and train quickly when transferred to new environments. Further, normative institutions will 537 be most effective in achieving collective alignment at scale when agent groups are large, as they provide structured guidance and shared norms that simplify decision-making and coordination across 538 many individuals. This reduces the complexity of aligning diverse actions and fosters widespread cooperation, even in expansive groups.

540 REFERENCES

553

561

562

565

566

567

570

577

578

579

584

585

586

542	John P. Agapiou, Alexander Sasha Vezhnevets, Edgar A. Duéñez-Guzmán, Jayd Matyas, Yiran Mao,
543	Peter Sunehag, Raphael Köster, Udari Madhushani, Kavya Kopparapu, Ramona Comanescu,
544	DJ Strouse, Michael B. Johanson, Sukhdeep Singh, Julia Haas, Igor Mordatch, Dean Mobbs,
545	and Joel Z. Leibo. Melting Pot 2.0, 2023.

- Abien Fred Agarap. Deep learning using rectified linear units (relu), 2019. URL https://arxiv.org/abs/1803.08375.
- Amritha Menon Anavankot, Stephen Cranefield, and Bastin Tony Roy Savarimuthu. Nemas: Norm
 entrepreneurship in multi-agent systems. *Systems*, 12(6):187, 2024.
- Tina Balke. A taxonomy for ensuring institutional compliance in utility computing. Schloss Dagstuhl-Leibniz Zentrum für Informatik, 2009.
- Tina Balke and Daniel Villatoro. Operationalization of the sanctioning process in utilitarian artificial
 societies. In *International Workshop on Coordination, Organizations, Institutions, and Norms in Agent Systems*, pp. 167–185. Springer, 2011.
- 557 Cristina Bicchieri. The grammar of society: the nature and dynamics of social norms. Cambridge
 558 University Press, 2005. URL https://api.semanticscholar.org/CorpusID:
 559 221193017.
 - Robert Boyd and Peter J Richerson. Punishment allows the evolution of cooperation (or anything else) in sizable groups. *Ethology and sociobiology*, 13(3):171–195, 1992.
- Amit Chopra, Leendert van der Torre, Harko Verhagen, and Serena Villata. *Handbook of normative multiagent systems*. College Publications, 2018.
 - Rosaria Conte and Cristiano Castelfranchi. Understanding the functions of norms in social groups through simulation. In *Artificial Societies*, pp. 225–238. Routledge, 2006.
- Robyn M. Dawes and David M. Messick. Social dilemmas. *International Journal of Psychology*, 2000.
- Christian Schroeder de Witt, Tarun Gupta, Denys Makoviichuk, Viktor Makoviychuk, Philip H. S.
 Torr, Mingfei Sun, and Shimon Whiteson. Is independent learning all you need in the starcraft multi-agent challenge? *arxiv:2011.09533*, 2020a.
- 574 Christian Schroeder de Witt, Tarun Gupta, Denys Makoviichuk, Viktor Makoviychuk, Philip H. S.
 575 Torr, Mingfei Sun, and Shimon Whiteson. Is independent learning all you need in the starcraft multi-agent challenge?, 2020b.
 - Yali Du, Joel Z Leibo, Usman Islam, Richard Willis, and Peter Sunehag. A review of cooperation in multi-agent learning. *arXiv preprint arXiv:2312.05162*, 2023.
- Christopher Frantz, Martin K Purvis, Mariusz Nowostawski, and Bastin Tony Roy Savarimuthu.
 nadico: A nested grammar of institutions. In *PRIMA 2013: Principles and Practice of Multi- Agent Systems: 16th International Conference, Dunedin, New Zealand, December 1-6, 2013. Proceedings 16*, pp. 429–436. Springer, 2013.
 - Julián García and Arne Traulsen. Evolution of coordinated punishment to enforce cooperation from an unbiased strategy space. *Journal of the Royal Society Interface*, 16(156):20190127, 2019.
- Michele J Gelfand, Sergey Gavrilets, and Nathan Nunn. Norm dynamics: Interdisciplinary perspectives on social norm emergence, persistence, and change. *Annual Review of Psychology*, 75(1): 341–378, 2024.
- Gillian K Hadfield and Barry R Weingast. What is law? a coordination model of the characteristics of legal order. *Journal of Legal Analysis*, 4(2):471–514, 2012.
- ⁵⁹³ Gillian K Hadfield and Barry R Weingast. Law without the state: legal attributes and the coordination of decentralized collective punishment. *Journal of Law and Courts*, 1(1):3–34, 2013.

594 595	Gillian Kereldena Hadfield. Rules for a flat world: Why humans invented law and how to reinvent it for a complex global economy. Oxford University Press, 2017.		
590 597 598	Dylan Hadfield-Menell and Gillian Hadfield. Incomplete contracting and AI alignment. <i>arxiv:1804.04268</i> , 2018.		
599 600	David Hales. Group reputation supports beneficent norms. <i>Journal of Artificial Societies and Social Simulation</i> , 5(4), 2002.		
601 602 603	Joseph Henrich. The secret of our success: How culture is driving human evolution, domesticating our species, and making us smarter. princeton University press, 2016.		
604 605	Matthew J Hoffmann. Self-organized criticality and norm avalanches. In <i>NORMAS</i> , pp. 117–125, 2005.		
606 607 608	Marco A. Janssen, Robert Holahan, Allen Lee, and Elinor Ostrom. Lab experiments for the study of social-ecological systems. <i>Science</i> , 2010.		
609 610 611 612	Natasha Jaques, Angeliki Lazaridou, Edward Hughes, Caglar Gulcehre, Pedro Ortega, DJ Strouse, Joel Z Leibo, and Nando De Freitas. Social influence as intrinsic motivation for multi-agent deep reinforcement learning. In <i>International conference on machine learning</i> , pp. 3040–3049. PMLR, 2019.		
613 614 615 616	Raphael Köster, Kevin R McKee, Richard Everett, Laura Weidinger, William S Isaac, Edward Hughes, Edgar A Duéñez-Guzmán, Thore Graepel, Matthew Botvinick, and Joel Z Leibo. Model-free conventions in multi-agent reinforcement learning with heterogeneous preferences. <i>arXiv</i> preprint arXiv:2010.09054, 2020.		
618 619 620	Raphael Köster, Dylan Hadfield-Menell, Richard Everett, Laura Weidinger, Gillian K Hadfield, and Joel Z Leibo. Spurious normativity enhances learning of compliance and enforcement behavior in artificial agents. <i>Proceedings of the National Academy of Sciences</i> , 119(3):e2106028118, 2022.		
621 622 623 624	Raphael Köster, Kevin R. McKee, Richard Everett, Laura Weidinger, William S. Isaac, Edward Hughes, Edgar A. Duéñez-Guzmán, Thore Graepel, Matthew Botvinick, and Joel Z. Leibo. Model-free conventions in multi-agent reinforcement learning with heterogeneous preferences. <i>arxiv:2010.09054</i> , 2020.		
625 626 627 628	Joel Z. Leibo, Julien Perolat, Edward Hughes, Steven Wheelwright, Adam H. Marblestone, Edgar Duéñez-Guzmán, Peter Sunehag, Iain Dunning, and Thore Graepel. Malthusian reinforcement learning. arxiv:1812.07019, 2019.		
629 630 631	Joel Z. Leibo, Edgar Dué nez Guzmán, Alexander Sasha Vezhnevets, John P. Agapiou, Peter Sune- hag, Raphael Koster, Jayd Matyas, Charles Beattie, Igor Mordatch, and Thore Graepel. Scalable evaluation of multi-agent reinforcement learning with melting pot. PMLR, 2021.		
632 633 634	Eric Liang, Richard Liaw, Philipp Moritz, Robert Nishihara, Roy Fox, Ken Goldberg, Joseph E. Gonzalez, Michael I. Jordan, and Ion Stoica. Rllib: Abstractions for distributed reinforcement learning, 2018. URL https://arxiv.org/abs/1712.09381.		
635 636 637 638	Moamin A Mahmoud, Mohd Sharifuddin Ahmad, Mohd Zaliman Mohd Yusoff, and Aida Mustapha. A review of norms and normative multiagent systems. <i>The Scientific World Journal</i> , 2014(1): 684587, 2014.		
639 640 641	Kevin R McKee, Ian Gemp, Brian McWilliams, Edgar A Duéñez-Guzmán, Edward Hughes, and Joel Z Leibo. Social diversity and social preferences in mixed-motive reinforcement learning. <i>arXiv preprint arXiv:2002.02325</i> , 2020.		
642 643 644 645	Kevin R McKee, Edward Hughes, Tina O Zhu, Martin J Chadwick, Raphael Koster, Antonio Garcia Castaneda, Charlie Beattie, Thore Graepel, Matt Botvinick, and Joel Z Leibo. A multi-agent reinforcement learning model of reputation and cooperation in human groups. <i>arXiv preprint arXiv:2103.04982</i> , 2021.		
646 647	Andreasa Morris-Martin, Marina De Vos, and Julian Padget. Norm emergence in multiagent sys-		

tems: a viewpoint paper. Autonomous Agents and Multi-Agent Systems, 33:706–749, 2019.

648 649 650	Julien Perolat, Joel Z. Leibo, Vinicius Zambaldi, Charles Beattie, Karl Tuyls, and Thore Grae- pel. A multi-agent reinforcement learning model of common-pool resource appropriation. <i>arxiv</i> :1707.06600, 2017.
652 653	Alexander Peysakhovich and Adam Lerer. Prosocial learning agents solve generalized stag hunts better than selfish ones. <i>arXiv preprint arXiv:1709.02865</i> , 2017.
654 655	Peter J Richerson and Robert Boyd. <i>Not by genes alone: How culture transformed human evolution</i> . University of Chicago press, 2008.
657 658 659 660	Bastin Tony Roy Savarimuthu, Stephen Cranefield, Maryam Purvis, and Martin Purvis. Role model based mechanism for norm emergence in artificial agent societies. In <i>International Workshop on Coordination, Organizations, Institutions, and Norms in Agent Systems</i> , pp. 203–217. Springer, 2007.
661 662 663 664	Bastin Tony Roy Savarimuthu, Maryam Purvis, Martin Purvis, and Stephen Cranefield. Social norm emergence in virtual agent societies. In <i>Declarative Agent Languages and Technologies VI: 6th International Workshop, DALT 2008, Estoril, Portugal, May 12, 2008, Revised Selected and Invited Papers 6</i> , pp. 18–28. Springer, 2009.
665 666 667	Bastin Tony Roy Savarimuthu, Surangika Ranathunga, and Stephen Cranefield. Harnessing the power of llms for normative reasoning in mass. <i>arXiv preprint arXiv:2403.16524</i> , 2024.
668 669	John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms, 2017. URL https://arxiv.org/abs/1707.06347.
670 671 672 673	Onkur Sen and Sandip Sen. Effects of social network topology and options on norm emergence. In <i>International Workshop on Coordination, Organizations, Institutions, and Norms in Agent Systems</i> , pp. 211–222. Springer, 2009.
674 675 676	Eugene Vinitsky, Raphael Köster, John P Agapiou, Edgar A Duéñez-Guzmán, Alexander S Vezhnevets, and Joel Z Leibo. A learning agent that acquires social norms from public sanctions in decentralized multi-agent settings. <i>Collective Intelligence</i> , 2(2):26339137231162025, 2023.
677 678 679 680 681 682 683 684 685 686 687 688 689 690 691 692 693 694 695 694 695 696 697 698 699	Yu Zhang and Jason Leezer. Emergence of social norms in complex networks. In 2009 International Conference on Computational Science and Engineering, volume 4, pp. 549–555. IEEE, 2009.

701

702 **DESCRIPTION OF ENVIRONMENT CONDITIONS** А 703

COMMONS HARVEST A.1

704

705 706

708

709

Vanilla: This environment retains the basic tagging (sanctioning) mechanism introduced by (Perolat et al., 2017). Tagging is costless for the sanctioning agent (other than opportunity cost) and removes the target agent from the environment for 25 steps. The tagged agent thus loses the opportunity to collect apples and the tagging agents benefits (if at all) from removing a competitor. This version 710 establishes a baseline for agent behaviors and group performance in the absence of norms, when 711 sanctioning can only generate private benefits for the sanctioning agent.

712 Hidden Rule SMG: This version incorporates pro-social rewards for sanctioning in accordance with 713 a norm. The norm prescribes which zone it is acceptable for agents to harvest from at a given point 714 in time. Initially, the acceptable zone is the one with the highest minimum apple count across its two 715 patches. This prescription changes dynamically: when one of the patches in the acceptable zone falls 716 below a threshold (four in our case) apples, the norm shifts to a zone with the highest minimum count 717 that meets the threshold. If no such zone exists, no zone is prescribed, and harvesting is prohibited 718 until regeneration occurs. This dynamic norm evolution ensures that acceptable behaviors adapt 719 to resource availability. We call this the 'hidden rule' condition because the norms/rewards for 720 sanctioning are generated by the environment and can only be discovered through sanctioning.

721 The tagging mechanism from the vanilla version is modified to enforce this normative structure. As 722 before, tagging costs the tagging agent -10 and tagged agents are removed from the environment 723 for 25 timesteps, losing harvesting opportunities. However, in addition to the private benefits to 724 tagging experienced in the vanilla version, if the target of sanctioning violated the current norm by 725 harvesting from a zone other than the one prescribed by the hidden reward structure, the tagging agent receives a reward of +20, resulting in a net reward of 10. This incentivizes agents to learn how 726 to correctly sanction, which is a dynamic problem. 727

728 Altared SMG: This condition implements the 729 same rewards for sanctioning and the same evo-730 lution of norms as the hidden rule condition but 731 also includes an altar - an observable classification institution incorporated at three distinct lo-732 cations within the environment, as shown in the 733 Figure 6. It serves as an environmental feature 734 encoding the currently prescribed norm. When 735 agents visit an altar location, they receive an ob-736 servation of its color, which matches the color 737 of the zone from which it is currently accept-738 able to collect apples. If no zone is acceptable 739 (i.e., all zones have insufficient resources), the 740 altar displays a yellow fire symbol. The sanc-741 tioning mechanism remains consistent with the 742 Hidden Rule SMG setup: agents earn rewards for tagging agents who collect from any zone 743 other than the one prescribed by the altar (or 744



Figure 6: Altared Commons Harvest: Altars display the color of the zone from which is it currently acceptable to harvest. Three zones: red (top two zones, 6 apples each), blue (middle two zones, 10 apples each), and green (bottom two zones, 13 apples each). Altar displays green: indicating that the bottom zone is acceptable for harvesting. The altar displays yellow fire when harvesting from all zones is prohibited.

from any zone when all zones are prohibited.) Unlike the hidden rule condition, where agents can 745 learn norms only from sanctioning and being sanctioned, in the Altared SMG the agents can also 746 learn to recognize the altar, to visit the altar, and to map its observations to appropriate sanctioning 747 behavior. 748

749

751

750 A.2 ALLELOPATHIC HARVEST

752 Vanilla: This environment simplifies the two layer zapping (sanctioning) mechanism introduced 753 by (Köster et al., 2020). Sanctioning is costless for the sanctioning agent (other than opportunity cost) and incurs a penalty of -10 for sanctioned agent. This version establishes a baseline for agent 754 behaviors and group performance in the absence of norms, when sanctioning can only generate 755 private benefits for the sanctioning agent.

756 Hidden Rule SMG: As before, this version incorporates pro-social rewards for sanctioning in accor-757 dance with a norm. The norm prescribes which monoculture is desired and planting the berry of that 758 color is acceptable action. This prescription changes dynamically: in the episode of 2000 steps, we 759 change the norm randomly every 100 steps for first 1000 steps and then 3-5 times at random interval 760 (minimum gap of 200 steps between change). This is to ensure that agents get enough experience for each color, while also mimicking real world processes such as regulatory updates. The zapping 761 mechanism from the vanilla version is modified to enforce this normative structure. As before, zap-762 ping costs the zapping agent -10 and tagged agents incurs cost of -10. However, in addition to the 763 private benefits to tagging experienced in the vanilla version, if the target of sanctioning violated the 764 current norm by planting the berry with color other than the one prescribed by the hidden reward 765 structure, the tagging agent receives a reward of +20, resulting in a net reward of 10. 766

Altared SMG: This condition implements the 767 same rewards for sanctioning and the same evo-768 lution of norms as the hidden rule condition but 769 also includes an altar - an observable classifica-770 tion institution incorporated at the center of the 771 environment, as shown in the Figure 7. When 772 agents visit an altar location, they receive an ob-773 servation of its color, which matches the desired 774 berry color to be planted. This will include 775 sanctioning gray agents too, thereby helping 776 towards solving fre-rider problem. The sanc-777 tioning mechanism remains consistent with the Hidden Rule SMG setup: agents earn rewards 778 for sanctioning agents who plant berry other 779 than the one prescribed by the altar. 780





B EXTENDED RELATED WORK

781 782

783 784

785

786

787

There is a vast body of literature addressing various aspects related to our agenda, reflecting extensive efforts across multiple directions. To provide a clear understanding of the existing work, we categorize the related efforts into several key topics, which are discussed in detail below.

Learning to cooperate in multi-agent systems. The problem of cooperation—how to design envi-788 ronments and algorithms that align agents' behavior towards higher collective welfare-has received 789 increasing attention in the multi-agent literature (Du et al., 2023). Common approaches include de-790 signing agents that have other regarding preferences through intrinsic rewards that promote collec-791 tive welfare Peysakhovich & Lerer (2017) or acting altruistically toward others (McKee et al., 2020). 792 Other methods use social influence Jaques et al. (2019) as an underlying mechanism; although these 793 approaches are generally designed for coordination problems rather than cooperation. It is worth 794 noting that in multi-agent reinforcement learning (MARL), the distinction between coordination 795 and cooperation challenges is often blurred. Recently, norms have emerged as another set of mechanisms in MARL specifically designed to address cooperation. These methods draw on the extensive 796 literature regarding the evolution of cooperation in human societies (Boyd & Richerson, 1992). By 797 extending Markov decision processes (MDPs) with sanctions (Vinitsky et al., 2023), agent societies 798 can support third-party punishment, and these methods have shown promise experimentally in fos-799 tering cooperation (Köster et al., 2022). However, most techniques still rely on direct modifications 800 of agent behavior through intrinsic rewards. These intrinsic reward depend on mechanisms such as, 801 mimicking others' punishment behaviors (Vinitsky et al., 2023) and developing positive reputations 802 (McKee et al., 2021). In contrast to these set of techniques, our method follows the key insight that 803 human societies did not learn to be cooperative just through exploration and individual behaviour 804 change. Rather, cooperation follows as a second-order effect once societies learn to coordinate their 805 peer sanctions through social structures, such as informal norms and formal institutions (Richerson 806 & Boyd, 2008; Henrich, 2016). Our work focuses on a particular manifestation of these structures, 807 namely, classification institutions, that announce right and wrong behaviours around which agents can voluntarily coordinate their sanctioning behaviour (Hadfield & Weingast, 2012). More impor-808 tantly, compared to previous work in MARL, we shift the focus from individual learning to learning about social structures. Specifically, our work uses standard MARL methods to give agents the

ability to recognize features of classification institutions (the *altars*) that represent the norms of a population.

Norms and Institutions in multi-agent systems. There is an extensive body of literature on norms 813 in multiagent systems (MAS), with frameworks addressing various stages of a norm life cycle, in-814 cluding norm emergence, transmission, enforcement, and internalization within artificial agent soci-815 eties (c.f Mahmoud et al. (2014); Chopra et al. (2018) for MAS literature and Gelfand et al. (2024) 816 for an interdisciplinary review). In MAS, institutions typically represent norms using formal declar-817 ative languages, similar to logical specifications. For instance, the nADICO framework Frantz et al. 818 (2013) provides a grammar for representing norms through institutional statements, and agents learn 819 the content and enforcement of these norms by observing the behavior of others. In contrast to these 820 symbolic norm representation methods, our approach uses a visual representation, eliminating the need for extensive handcrafted specifications in a formal language. Within the context of learning, 821 another key difference lies in the focus of our work on learning enforcement behaviour rather than 822 learning norm compliance Savarimuthu et al. (2024). In institutionalized MAS, norm enforcement 823 through sanctions against violating agents is often a centralized process. Even when sanctions are 824 imposed by third parties, the enforcer is typically a specially designed agent with dedicated mon-825 itoring roles Balke (2009); Balke & Villatoro (2011). In our approach, however, enforcement is 826 entirely decentralized by making it part of the learned behaviour of each and every agent. In a 827 related work, Garcia and Traulsen García & Traulsen (2019) analyze the effects of different pool 828 punishment institutions, specifically pro-social and anti-social centralized institutions, where mem-829 bers can contribute to a coordinated punishment scheme. The model finds that public visibility of 830 pro-social institutions is essential for the stability of cooperative strategies, as agents can condition 831 their behavior based on the presence and visibility of these institutions. In comparison, our work do not deal with the question of establishment of the institution as part of the agent strategy, however, 832 we analyze the impact of different types of institutions as well as visibility of institutions. 833

834 Norm creation and emergence. Embedding normative behavior in agents is commonly referred 835 to as norm creation in the multiagent systems (MAS) literature Chopra et al. (2018). Previous 836 approaches often treated this as an offline process, where agents were pre-programmed to follow 837 specific norms, such as those related to property rights Conte & Castelfranchi (2006) or reputation (Hales, 2002). More recent approaches have introduced models of norm creation through specialized 838 agents known as norm entrepreneurs Savarimuthu et al. (2007); Anavankot et al. (2024). However, 839 norm creation through specialized agents raises additional questions, such as how and why norms 840 are accepted and transmitted in a population Hoffmann (2005). This introduces the factor of network 841 topology that determines agent interactions, and adds another layer of complexity to the process Sen 842 & Sen (2009). Whereas norm creation is analyzed at a micro-level, norm emergence is studied as a 843 macro phenomenon in artificial societies, often driven by a threshold effect—if a certain proportion 844 of the population adheres to a behavior (descriptive norm) and enforces or expects the enforcement 845 of that behaviour (social norm), that behavior can become widespread Morris-Martin et al. (2019). 846 In this context, various models analyze the impact of independent variables at the micro-level, such 847 as the cost of enforcement Savarimuthu et al. (2009), and environmental factors, such as network 848 topology Zhang & Leezer (2009), on norm emergence.

849 850

C TRAINING

851 852

In our experiments, we employed RLlib Liang et al. (2018) for training, utilizing the Proximal Pol-853 icy Optimization (PPO) algorithm Schulman et al. (2017) to optimize agent behaviours. The agent 854 architecture consisted of fully connected layers with hidden sizes of 64 and 256, using ReLU activa-855 tions Agarap (2019). The environment was based on DeepMind's Melting Pot library Agapiou et al. 856 (2023), with custom configurations designed to match the specific task objectives. We performed 857 hyper-parameter tuning on network sizes and learning rate, train batch size and CNN filters using 858 grid search on 200 configurations. Table 1 contains a list of parameters used during training. All 859 experiments were run on a single GPU node with one A40 GPU and 32 CPUs. For this work, we 860 consider a simple, practical and scalable approach to MARL is Independent Proximal Policy Opti-861 mization (IPPO), a decentralized method where each agent optimizes its policy independently using a variant of Proximal Policy Optimization (PPO). IPPO is well-suited for learning in multi-agent 862 systems due to its ability to stabilize learning through clipped policy updates while maintaining 863 scalability. Additionally, IPPO can be further enriched through inputs such as agent indication,



Figure 8: Agent architecture and interaction with environment.

which assigns unique identifiers to agents, and agent roles, which guide agents to adopt distinct strategies. These enhancements facilitate the learning of heterogeneous and independent policies, increasing the framework's flexibility and adaptability to diverse multi-agent dynamics. These capabilities make IPPO the algorithm of choice for our investigation, allowing us to effectively study compliance and coordination in dynamic, multi-agent settings.

898 C.1 AGENT ARCHITECTURE

899 The key objective of this work is to assess the implications of normative infrastructure on the align-900 ment and cooperation capabilities of the agents. We are indeed proposing to shift the focus away 901 from building complicated agent architectures in order embed norms and value in them and point the 902 focus towards building normative infrastructure. Given this, we chose to conduct our experiments 903 using a very simple shared-parameter architecture of our agents consisting of convolutions layers 904 followed by fully connected networks. We tested both with and without GRU units and did not 905 find significant difference in performance. Our architecture closely follows independent learning 906 PPO de Witt et al. (2020b) agents, where the actor and the critic network are shared between agents. For heterogeneity in independently learning agents, we include both the agent indication and their 907 pre-defined roles in the environment as extra input. We believe that our proposal is agnostic to the 908 agent architecture and testing with more sophisticated agents is inteded as future work. Training 909 details are available in Appendix C. 910

911 912

913

890 891 892

893

894

895

896 897

D ADDITIONAL RESULTS

914 D.1 ALLELOPATHIC HARVEST WITH FIXED ALTAR 915

Here, we consider an ablation of our Altared Allelopathic Harvest environment where the altar
 remains stationary throughout the episode (across all episodes) instead of being dynamic. The environment details remain the same as mentioned in Section E.1.1 apart from the changes discussed

918	Parameter	Value	
919		Resources	
920	Number of Rollout Workers	30	
921	Number of GPUs	1	
922		Training	
923	Seeds Used	12345, 67890, 54321, 98765, 20242	
924	Rollout Fragment Length	100	
925	Train Batch Size	32,000	
926	SGD Minibatch Size	4,096	
927	Number of SGD Iterations	30	
928	Disable Observation Preprocessing	True	
929	Use New RL Modules	False	
930	Use New Learner API	False	
931	Framework	torch	
022	Agent Model		
022	Fully Connected Hidden Layers	(64, 64)	
933	Post-FC Hidden Layer	(256)	
934	CNN Activation	ReLU	
935	FC Activation	ReLU	
936	Post-FC Activation	ReLU	
937	LSTM Use Previous Action	True	
938	LSTM Use Previous Reward	False	
939	LSTM Cell Size	256	
940	Experiment Trials		
941	Stopping Criteria	10,000 training iterations	
942	Number of Checkpoints	30	
943	Checkpoint Interval	50	
944	·	·,	

Table 1: Training and Hyper-parameter Configuration

in this section. We choose the altar to display red colored berry, making red monoculture the desirable outcome. We test three distinct conditions (similar to the versions presented in Section E.1.1) to explore the effects of different sanctioning mechanisms on agent's enforcement and compliance behavior and their ability to achieve a monoculture. We discuss these below:



hing agents. Please note that we have adjusted the

Figure 9: Rewards of training agents. Please note that we have adjusted the reward curves to only include reward for berry consumption and penalty for getting sanctioned while removing any effect of reward and cost received due to sanctioning rules

Results Figure 9 shows agent performance in maximizing welfare, measured as the sum of rewards across agents and averaged over episodes. Welfare is maximized when agents align their planting and sanctioning behaviors to achieve a monoculture of one berry color, as faster ripening occurs with





Figure 13: Monoculture and Agents' status at halfway of an episode for trained agents

E ENVIRONMENT DETAILS

1053 E.1 Allelopathic Harvest

Background and Setup The 'Allelopathic Harvest' environment 1055 (Agapiou et al., 2023; Köster et al., 2020) is a mixed-motive game 1056 which poses both the coordination and the free-rider problem, mak-1057 ing it challenging for agents to reach a welfare maximizing out-1058 come. It features a map containing a total of 348 berries of three 1059 different colors (116 of each red, blue and green) and sixteen agents that can plant and consume berries. Each agent has an intrinsic 1061 preference for a specific color berry. Out of 16 agents, 8 prefer 1062 red berries and other 8 prefer green berries by default. Agents get 1063 reward for consuming any ripened berry (+1) but receive higher re-1064 ward for consuming their preferred color berry (+2). Agents can also plant berries of specific color but that does not generate any reward or cost and hence agents have no direct incentive to plant, 1066 leading to a free-rider problem. After planting a berry, agent's color 1067 changes to the color of the planted berry. However, after eating a 1068 ripened berry, their color is stochastically reset to gray The agents 1069



Figure 12: > 95% red monoculture.

can only consume ripened berries and the berry ripening rate is directly proportional to the fraction of the amount of berry of that color. Agents also have a zapping action which fires a white beam that they can use to tag other agents. When an agent is zapped (target), it receives a penalty of -10. While zapping, source agents don't receive any reward or penalty by default (but this may be changed in different versions of the environment below). An episode of this environment lasts 2000 timesteps. More details about it can be found in Agapiou et al. (2023).

1075

1049 1050 1051

1052

1054

1076 E.1.1 TRAINING ENVIRONMENTS 1077

Altar In this version, we introduce a **dynamic** 'altar' (Fig. 7) in the map – a visual observation $(3 \times 3 \text{ subgrid})$ in the center of the map that displays berries of a specific color whose monoculture is desired. Agents have an augmented observation space that includes a memory slot, which starts as

1082





1093

1094

1095

1096

(a) Altar Prescribed Berry: Red

(b) Altar Prescribed Berry: Blue (c) Altar Prescribed Berry: Green

Figure 14: Illustration of the Altared Allelopathic Harvest environment across different timesteps in an episode where the altar is present in the center of the environment depicting colored berries. The altar prescribed color changes periodically during an episode.

empty. When an agent enters a tile that is part of the altar, their memory slot updates to altar obser-1099 vation (Fig. 8). The altar is not stationary and its color changes dynamically. The altar color at the 1100 start of the episode is set to be randomly among red, blue and green. For the first 1000 timesteps, the 1101 altar changes color in a fixed manner at every 100 timesteps. However, after that, the color changes 1102 partially randomly in the following manner: nextUpdateStep = previousUpdateStep 1103 + 160 + random(1, 100) where nextUpdateStep denotes the timestep (in future) when the color needs to be changed. At each update, the next color is chosen such that it is not the same 1104 as previous color. More specifically, the next color is sampled randomly from the two remaining 1105 colors with equal probability for each. 1106

1107 Further, the presence of the altar also influences the reward dynamics associated with zapping (tag-1108 ging). Specifically, if a source agent zaps a target agent of the same color displayed on the altar at 1109 that moment, both the source and target agents receive a penalty of -10 points. If the source agent 1110 zaps a target of any other color, the source agent receives a net reward of +10 points, while the target still receives a penalty of -10 points. Thus, the altar essentially only influences the reward (penalty) 1111 received by source (zapping) agent. There is no reward (penalty) if the zapping beam doesn't hit any 1112 agent. 1113

1114 We refer to this environment as Altared Allelopathic Harvest and provided an overview 1115 of it in Figure 14.

1116

1117 Hidden Rules In the 'Hidden Rules' variant, the physical altar is removed from the environment, 1118 while all other dynamics remain the same as in the Altared Allelopathic Harvest ver-1119 sion. Thus, the altar's influence persists only in the form of controlling tagging (zapping) related 1120 rewards and penalties without visual presence. If a source agent zaps another agent that is the same 1121 color as prescribed by the (hidden) altar, both the source and target agents receive a penalty of -10 1122 points. If the source agent zaps an agent of any other color, the source agent receives a net reward 1123 of +10 points. We note that here also, even thought the altar is hidden, its color updates in the same 1124 manner as described in the Altared variant.

1125

1126 **Vanilla (Free Sanctioning)** In this condition, there is no altar or hidden rule in the environment. 1127 Agents can freely zap other agents, with the target agent receiving a penalty of -10 points. The 1128 source agent, however, does not receive any reward or penalty for zapping.

1129 1130

1132

1131 **E.1.2** EVALUATION ENVIRONMENTS

In this section, we describe some extensions of the Altared Allelopathic Harvest envi-1133 ronment which we used to evaluation of the agents trained in environments described earlier.

Limited Altar In this variant, the altar vanishes from the environment and reappears periodically. At every nextUpdateStep, we alternate between either removing the altar from the environment (without changing its color so the agents still receive reward or penalty upon zapping other agents) or updating the color of the altar and reinstantiating it visually in the environment. For example, for the first 1000 timesteps, the altar disappears at 100, 300, 500, 700, 900 step while it updates color and reappears at 200, 400, 600, 800, 1000 step. When the altar is not present, the agents see blank (empty) cells on its place and don't receive any observation when stepping into it.

Thus, the altar is only visible and accessible for a limited number of steps. The objective of this environment is to test whether agents learn to optimize their visits to the altar to maximize updated knowledge within a restricted window. Fig. 15 provides an overview of this environment.



Figure 15: Illustration of 'Limited Altar' version of the Altared Allelopathic Harvest environment.

1161 1162

1151 1152 1153

1163 Conflicting Altars In this environment, we introduce two more secondary 'altars' of size 2 × 2 1164 and place them at the bottom-right and top-left regions of the map while the primary altar remains 1165 in the center. All three altars change their color periodically (at the same nextUpdateStep) but 1166 the secondary altars always show a color contradicting the primary (center) altar's prescribed color. 1167 However, the underlying reward dynamics associated with zapping only depends on the primary 1168 altar and changes when it updates. A key point here is when an agent visits a secondary altar, its 1169 observation gets updated wrongly.

The objective of this environment is to test whether the agents can figure out which altar (normative institution) is the "correct" one by visiting them and interpreting the signals received and then only choosing to visit the correct one at convergence. Fig. 16 provides an overview of this environment.

1173

1174 E.2 COMMONS HARVEST

Background The "Commons Harvest" environment is inspired by the 'Commons Harvest' substrate in Melting Pot (Agapiou et al., 2023), which itself draws from prior work on multi-agent
reinforcement learning for common-pool resource appropriation (Perolat et al., 2017). In this environment, agents aim to collect apples scattered across six distinct patches. Each patch consists of
multiple apple cells, with each cell having at least one neighboring apple.

Agents receive a reward of 1 for every apple consumed. Apples regenerate with a per-step probability that depends on the number of neighboring apples within an Euclidean distance of 2. Exact details about the regrowth probability can be found in Agapiou et al. (2023). Specifically, if there are no apples in the vicinity, the probability of regrowth is zero. Consequently, patches can be **permanently depleted** if all apples in a patch are harvested, requiring agents to exercise caution and avoid overharvesting. If agents exhaust a patch, it will not recover, and sustaining apple regeneration demands collective restraint among the agents. This dynamic leads to a *social dilemma*, akin to the tragedy of the commons, where individual incentives clash with the group's long-term interest.



Figure 16: Dynamic Conflicting Institutions

Figure 17: Illustration of 'Conflicting Altars' environment where the secondary altars at the bottom right and top left region of the environment always show a color different from the primary altar (in the center) to distract the agents.

Agents face a strong incentive to consume the last remaining apple to maximize individual gain, potentially at the cost of losing that patch permanently.



1223 (a) Single apple left in one of the patches with a blue (b) Blue agent eats the last remaining apple leading to 1224 agent standing next to it.

the patch being lost permanently.

Figure 18: **Commons Harvest**: Illustration of how a patch can be lost irrevocably if the last remaining apple is eaten by any of the agents.

1227 1228

1225

1226

1202 1203

1204

1205

1206 1207 1208

1209

1229 **Common Setup** The environment features seven agents, six apple patches, and each episode runs 1230 for 5000 timesteps. Agents are initially colored gray. The six apple patches are grouped into three 1231 zones ---- red (zone 1), green (zone 2), and blue (zone 3) ---- with each zone containing two patches. When an agent eats an apple from a patch, its avatar color changes to the color of the corresponding 1232 zone, allowing agents to observe from which zone others have recently eaten. 1233

1234 As in the related "Allelopathic Harvest" environment, agents can also tag each other with a beam. If 1235 an agent is tagged, it is removed from the environment for 25 steps. No direct reward or punishment 1236 is received for tagging or being tagged, but there are indirect consequences: the tagged agent loses 1237 opportunities to collect apples during its timeout, while the tagging agent faces the opportunity cost of spending time on tagging rather than gathering apples. 1238

1239

Altared Version In this setting, in addition to the existing setup, we introduce an 'altar' in the 1240 map - a visual indicator located at three positions: bottom left, bottom right, and center of the map. 1241 The altar consists of a 2x2 grid and displays the color of the zone from which agents should ideally



Figure 19: A snapshot from the 'Altared Version' of the Commons Harvest environment. Each colored box denotes an apple patch where the color indicates the zone to which the patch belongs. Agents obtain the color of the zone from which they last ate. Altar (2×2 subgrid) is present in center, bottom left, and bottom right regions of the map. It is currently colored green indicating that apples should ideally be eaten from the green zone.

1270

1271

1272

consume apples. When an agent enters the altar cells, it observes the altar color, which changes
dynamically based on the number of apples remaining in each zone. A illustration of the different
zones and altar is shown in Fig. 19 and the dynamics of the norms is illustrated in Figure 21.

The altar's color is initially set to the color of one of the zones whose patches' minimum apple count is maximum overall. It remains that color until one of the patches in the associated zone has less than 4 apples in which case its color is set to the zone whose patches' minimum apple count is maximum overall and above 3 at that moment. If no zone has both patches with more than 3 apples, the altar displays a yellow fire symbol, signaling that agents should refrain from consuming apples from any patch to avoid the risk of permanently losing them. The altar remains in the fire state until apple regeneration occurs.

1283 When an agent tags another agent, the tagged agent is removed from the environment for 25 1284 timesteps, similar to the original setup. However, the tagging (source) agent incurs a penalty of 1285 -10 if the tagged agent's color matches the altar color, indicating that the target agent was follow-1286 ing the altar's guidance. Conversely, if the tagged (target) agent has a different color (excluding 1287 gray), the tagging agent receives a reward of +10. We show an illustration of the altar-based reward 1288 dynamics in Fig. 22a. When the altar is displaying fire, tagging any non-gray agent results in a +101289 reward.

1290

1291

Hidden Rules In the 'Hidden Rules' variant, the physical altar is removed from the environment,
while the reward and penalty mechanisms remain the same as in "Altared Version". Thus, the altar's influence persists only in the form of controlling tagging (zapping) related rewards and penalties without actual presence. We show an illustration of this in Fig. 22b. Since agents no longer observe the altar directly, they would have to infer the altar's state based solely on the rewards they receive.



1312 1313 1314 1315 1316 1317 1318 1319 (a) Eat from blue zone (b) Eat from red zone (c) Eat from green zone (d) All zones prohibited

Figure 21: Dynamics of ALTARED Commons Harvest



(a) Altared Version: The color of the altar is green and the grey agent fires a beam hitting both a green and red agent. The grey agent receives a penalty of -10 for zapping the green agent as its color is same as the altar but it receives a reward of +10 for zapping the red agent. Both the tagged agents are removed and reappear in the map after 25 timesteps.

(b) **Hideen Rules**: The altar is not present physically in the environment but its color is green. The grey agent fires a beam hitting both a green and red agent. It receives a penalty of -10 for zapping the green agent as its color is same as the virtual altar but it also receives a reward of +10 for zapping the red agent. Both the tagged agents are removed for 25 timesteps.

Figure 22: Illustration of the reward dynamics based on the altar in different versions of the Common Harvest environment.

1342 1343

1334

1335

1336

1337

1338

1339

1320

1321

Vanilla In the 'Vanilla' variant, the altar is not present, and the rewards and penalties associated
with the altar are also eliminated. The agents receive rewards only for consuming apples, and receive
no other signal from the environment. A snapshot of this environment at the beginning is shown in
Fig. 20b.

- 1348
- 1349