

Silly rules enhance learning of compliance and enforcement behavior in artificial agents

Raphael Köster^{a,1}, Dylan Hadfield-Menell^{b,c}, Richard Everett^a, Laura Weidinger^a, Gillian K. Hadfield^{c,d,e,f}, and Joel Z. Leibo^{a,1}

^aDeepMind; ^bDepartment of Electrical Engineering and Computer Science, University of California Berkeley; ^cCenter for Human-Compatible AI; ^dSchwartz Reisman Institute for Technology and Society, University of Toronto; ^eVector Institute; ^fOpenAI

This manuscript was compiled on April 25, 2021

1 **How do societies learn and maintain social norms? Here we use**
2 **multi-agent reinforcement learning to investigate the learning dynam-**
3 **ics of enforcement and compliance behaviors. Artificial agents pop-**
4 **ulate a foraging environment and need to learn to avoid a poisonous**
5 **berry. Agents learn to avoid eating poisonous berries better when do-**
6 **ing so is taboo, meaning the behavior is punished by other agents.**
7 **The taboo helps overcome a credit-assignment problem in discov-**
8 **ering delayed health effects. By probing what individual agents**
9 **have learned, we demonstrate that normative behavior relies on a**
10 **sequence of learned skills. Learning rule compliance builds upon**
11 **prior learning of rule enforcement by other agents. Critically, intro-**
12 **ducing an additional taboo, which results in punishment for eating**
13 **a harmless berry, further improves overall returns. This “silly rule”**
14 **counterintuitively has a positive effect because it gives agents more**
15 **practice in learning rule enforcement. Our results highlight the ben-**
16 **efit of employing a computational model focused on learning to im-**
17 **plement complex actions.**

Multi-agent reinforcement learning | Norms | Third-party punishment

1 **O**ne of the central attributes that differentiates human from
2 other animal societies and accounts for the enormous
3 gains of human ultra-sociality (1) is the presence of third-
4 party enforced norms (2–4). Many of these norms generate
5 direct benefits for individual and group well-being: norms
6 that prescribe reciprocity, fair sharing of rewards, or non-
7 interference with property properly claimed by another, for
8 example, can coordinate behavior and sustain incentives for
9 cooperation and investment. These are the norms that are the
10 primary focus of most research into the properties and origins
11 of human normativity (see (3) for a review.)

12 The normative landscape is also, however, populated by
13 many norms that appear essentially arbitrary: norms about
14 how and what we eat, how we greet each other, what clothes
15 and body decorations we wear, and what rituals we observe
16 (5, 6). People treat compliance with these norms as impor-
17 tant and punish violations, but, except for effects generated
18 by this socially-constructed salience, they have no direct or
19 first-order impact on welfare. Fessler et al. (5) call the process
20 by which patterns of behavior are imbued with moral senti-
21 ments that motivate sanctioning of violations of the pattern
22 *normative moralization*. They use as an example the norma-
23 tive moralization of handedness. Most people are naturally
24 right-handed but, particularly in societies with few special-
25 ized tools, whether someone is right- or left-handed generally
26 has no material consequences for others. Nonetheless, many
27 cultures treat using one’s right hand as a morally approved
28 category—denoting purity or politeness—and one’s left hand as
29 cause for opprobrium—revealing weakness or evil (7). Following
30 Hadfield-Menell et al. (8) we call such social norms *silly rules*.

31 The ubiquity of silly rules provides a puzzle for functionalist
32 accounts of norms (9); several explanations have been explored
33 so far. One kind of explanation posits that silly rules may exist
34 to serve as cheap signals of group membership and thus facili-
35 tate cooperation within the group (10). Another account holds
36 that silly rules are stable because, in any society, the survival
37 of each generation depends on the transmission from prior gen-
38 erations of a large amount of culture-specific and, importantly,
39 *causally opaque* knowledge (11). This includes everything
40 from which local plants produce edible versus poisonous fruit,
41 to how best to organize to resolve disputes between family
42 members. Most of the time individuals have no way of know-
43 ing which of the many rules they follow are critical for their
44 well-being. Thus silly rules may remain stable by virtue of
45 their incorporation into larger normative systems that also in-
46 clude important rules (1). Further support for this hypothesis
47 is found in the tendency of human children to over-imitate
48 adults, copying—and moralizing—even apparently irrelevant
49 aspects of adult behavior (12). The sheer abundance of silly
50 rules seems to require an account that grants the normative
51 moralization of seemingly irrelevant actions a more significant
52 role. It would seem that a society would do better to minimize
53 costly efforts to punish and conform with norms that produce
54 no material benefits, and so to economize on the number of
55 silly rules used as markers or retained as a by-product of the
56 cultural transmission of knowledge.

57 In this paper, we describe a new kind of functional ex-
58 planation for silly rules based on the dynamics of learning
59 in a society that lacks *a priori* knowledge of which of their
60 rules are truly important (causal opacity). Our explanation
61 relies on an essential asymmetry between the enforcement and
62 compliance aspects of normative behavior. In short, the skills
63 involved in third party norm enforcement readily transfer from
64 norm to norm, while the skills involved in compliance are
65 norm-specific. Thus adding a silly rule to a normative system
66 that already contains some number of deeply important rules
67 can be beneficial because the silly rule may provide greater op-
68 portunity to practice third party norm enforcement, a generic
69 skill. Improved norm enforcement by the group then makes
70 it easier for individuals to learn from experience the skills
71 necessary for norm compliance, such as how to prospectively
72 recognize and avoid specific taboos. Therefore, introducing a
73 silly rule may positively impact the learnability of compliance
74 behavior for all of a society’s rules, including those that truly
75 are important. The benefit of learning important rules faster

RK, RE, LW, JZL conducted the simulation and analysis. All authors contributed to the research design and writing of the paper.

The authors declare no conflicts of interest.

¹To whom correspondence should be addressed. E-mail: rkoster@google.com, jzl@google.com

76 can easily outweigh the dead-weight loss created by the silly
77 rule.

78 Silly rules can support the emergence and stability of a ben-
79 efiticial *normative social order* (13). In a normative social order,
80 group behavior is patterned on a classification scheme (called
81 a norm) that divides behaviors into approved and disapproved
82 (taboo) categories. Here, we employ a computational approach
83 to investigate the effects of silly rules on how well a normative
84 social order is learned. Our model consists of a multi-agent
85 reinforcement learning (RL) environment with eight artificial
86 agents, all simultaneously learning and interacting with one
87 another. Agents in our environment (Fig. 2) are faced with
88 learning a foraging task: learning to find and consume food
89 (“berries”). We assume that berries are relatively abundant
90 so there is no competition between agents and no common
91 pool resource problem. What makes the environment chal-
92 lenging is the presence of a poisonous berry which if eaten
93 will then reduce the value of an agent’s future consumption.
94 But importantly, the deleterious effect only triggers after a
95 significant delay. The delay introduces a credit-assignment
96 problem, meaning it is difficult for our agents to learn which
97 particular berry caused the negative effect and thus to learn
98 to avoid it. In this setting, a taboo on the poisonous berry—
99 however it may evolve—raises individual welfare. As Boyd et
100 al. (11) emphasize, this is a critical pathway by which culture
101 raises human well-being: through the transmission of cultural
102 practices, such as the avoidance of harmful foods, even when
103 agents lack direct causal awareness of why their practices are
104 beneficial (11, 14). The mechanisms behind such social learn-
105 ing in humans may be multiple: a psychological propensity
106 to conformity (15, 16), deliberate teaching practices (17), and
107 third-party punishment of failures to follow norms (18). Our
108 work focuses on the last of these. We show that agents are
109 able to sustain the transmission of a valuable taboo in order to
110 avoid a poisonous berry. For this, agents need to have learned
111 to recognize when another agent has violated a taboo and to
112 deliver a costly punishment to the violator.

113 Because our model allows us to separate the learning of
114 enforcement and compliance behaviors from the learning of
115 norm content itself, we designed an experiment in which norm
116 content was fixed in advance by the experimenter (which color
117 berries were taboo). By varying the content of the norms, we
118 can study the downstream effects on how the normative social
119 order (enforcement and compliance behavior) is learned. If a
120 player breaks a taboo they change color and become ‘marked’.
121 We assumed all agents have a form of mutual knowledge of
122 the rule in the sense that they may perceive violations of other
123 players via the marking. Note that players cannot directly
124 perceive their own marking, otherwise the self-marking would
125 trivially solve the credit assignment problem. If a player is
126 marked, other players can collect a reward for punishing them.
127 This creates an incentive for players to learn to punish rule
128 violations, and for players to not violate any rules. This reflects
129 situations in which there is a centralized scheme that labels
130 transgressive behavior, but the enforcement is decentralized.
131 For example, in medieval iceland the ‘law speaker’ would label
132 acts as unlawful. Individuals would be declared ‘outlaws’, and
133 others could take their property without repercussions (19).

134 Our environment is intended to capture the evolutionary
135 challenge humans faced in developing the behavioral and cog-
136 nitive repertoires of normativity and third-party punishment,

137 attributes that distinguish humans from other primates and
138 account for the gains humans enjoy from ultrasociality and
139 extraordinary gains from cooperation (2, 3, 20, 21). It is
140 the challenge of discovering and learning these behaviors—to
141 punish violators and to avoid punishment through behavior
142 modification—that animates our study. We demonstrate that
143 simple RL agents that lack any “mental” models of rules, vi-
144 olations, or punishment, will nonetheless learn to enforce the
145 rules. The enforcement subsequently enables agents to learn
146 to comply with a taboo against eating the poisonous berry.

147 As a precondition to our main analysis, we show that
148 individuals achieve higher overall welfare in a world where
149 eating the poisonous berry is taboo, relative to a world in
150 which there are no taboos so avoidance of the poisonous berry
151 must be learned only through individual experience. We show
152 that even with the cost of enforcement, overall group welfare is
153 higher with a norm than without. Thus, the normative social
154 order is valuable. We then show our main result: that the
155 value of a normative order is *higher* if the norms in this regime
156 include not only important rules—such as the rule against
157 eating poisonous berries—but also silly rules which make the
158 eating of a harmless berry taboo and bring about the same
159 third-party punishment. In our environment, agents learn
160 to enforce, and comply with, norms more quickly if the rule
161 system includes two taboos—one against eating the poisonous
162 berry, and one against eating a harmless berry.

163 Our results demonstrate a new account of the ubiquity
164 of rules that have no first-order impact on well-being. They
165 also provide a formalization of normativity in a computational
166 setting that we think will expand the tools available both
167 for understanding how human normativity operates, and how
168 artificial agents that are capable of participating in human
169 normative social orders might be built. In this sense, this
170 work is part of a research program that ultimately aims to
171 develop models capable of capturing distinctive features of
172 human intelligence such as the origin of institutions (22).

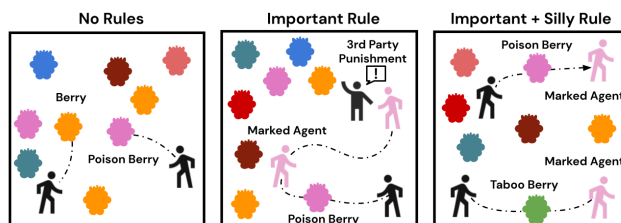


Fig. 1. Schematic overview of the experimental conditions. In the *no rules* condition, agents collect berries for reward. One berry-color is poisonous and after a time delay reduces reward obtained from consumed berries. Being poisoned is invisible to all players and a hard credit-assignment problem. In the *important rule* condition, eating the poisonous berry is a social taboo. When eaten, the player who ate the berry immediately gets marked, which is only visible to other agents. Other agents can collect a reward by punishing a marked agent. In the *important+silly rule* condition, the same taboo against the poisonous berry is in place, but additionally there are an identical taboo on a berry that is not poisonous. Therefore, there are two taboos for which agents can experience punishment by others. Our experiment sets out to study the effects of these different normative schemes.

Studying social norms with deep reinforcement learning. Computational simulations of populations and cultural development typically use an abstracted or idealized space to encode normative structure (23–26). Agents are usually mod-

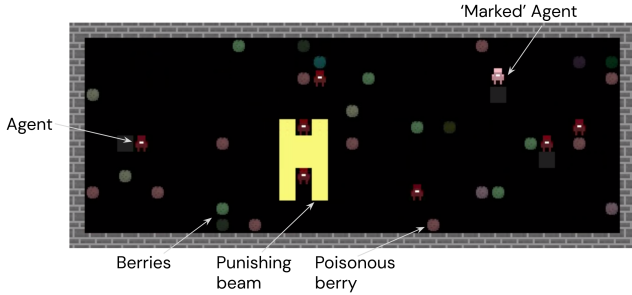


Fig. 2. Depiction of the environment. The agents inhabit a grid world. Agents earn reward for eating berries, which regrow probabilistic after being harvested. One type of berry is poisonous and if collected by an agent, it diminishes the agent’s ability to gather rewards from other berries, after a delay period. If an agent eats one of the poisonous berries in the *important rule* condition, the agent immediately gets “marked” and appears in a different color to the other agents. In the *important+silly rule* condition, one additional, non-poisonous, berry also triggers an agents’ marking. Agents are able to punish each other using a “punishing beam”, causing a loss to themselves and a large loss to the punished agent. If a “marked” agent is punished, the punishing agent receives a large reward.

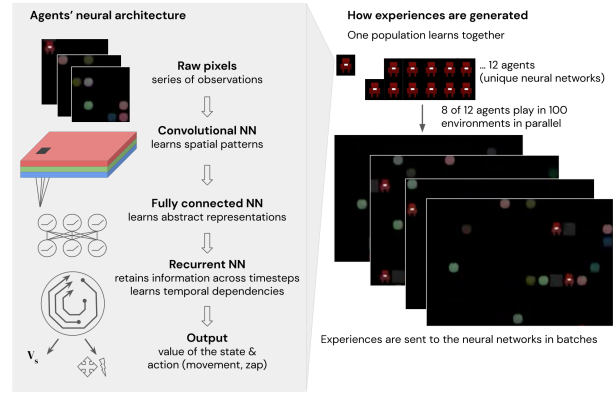


Fig. 3. Agent architecture and training procedure. Agents learn together in one population of 12 agents, 8 of which are selected to play in one episode in order to generate experiences (in multiple parallel environments). Each agent contains an independent neural network that receives a batch of its own experiences from these environments to update its neuronal weights. The inputs the agents receive are the raw pixels from their field of view. The network architecture of each agent consists of a convolutional neural network that learns to decompose the input into spatial patterns. This projects to fully connected layers that learn more abstract representations of game states, followed by a recurrent network that is able to retain and transform information over multiple timesteps. The output of neural network on each timestep is a prediction of the value of the current state and an action (movement or zap). Network weights are gradually adjusted to maximize long term cumulative reward.

177 eled either as choosing strategies within a game theoretic
 178 framework in which they are supplied with a set of available
 179 actions and associated payoffs, or as implementing behavioral
 180 rules in competition with other similarly-constructed agents.
 181 In these approaches, agents choose what to do (e.g. cooperate
 182 or defect), but the models cannot capture phenomena related
 183 to how they learn to implement their choice.

184 Here we apply a more generalized framework, a multi-agent
 185 RL approach that has been successfully used to study in-
 186 tertemporal (sequential) social dilemmas (27–37). Agents in
 187 this framework are artificial neural networks, which learn be-
 188 havioral policies (associating actions to states) and obtain
 189 rewards from an environment. State transitions in the environ-
 190 ment are generated by the actions of all agents combined.
 191 Agents inhabit a 2-D world in which they and other objects
 192 are located at coordinates in space. An agent’s action-space
 193 consists of moving up, down, left, right, rotating left and right
 194 and using a “punishing beam” directed at an adjacent agent
 195 (Fig. 3). Use of the punishing beam costs the punisher 20
 196 points and inflicts a cost of 35 on the punished agent. A vari-
 197 ety of ‘berries’ of different colours are distributed randomly
 198 throughout the world. An agent receives a reward of 4 if it
 199 navigates to a square with a berry, interpreted as “eating it”.
 200 Berries grow in sufficient abundance that there is no competi-
 201 tion between agents. Pink berries are poisonous: 100 timesteps
 202 after consumption they reduce reward gained by future berries
 203 to 1 point.

204 The environment can include a latent classification scheme
 205 (13) that designates some berries (colors) as “taboo”. This
 206 normative classification is implemented by inducing a change
 207 in the color of an agent (visible only to other agents) who
 208 consumes a taboo berry and changing the payoff associated
 209 with use of the punishing beam against such an agent. Pun-
 210 ishing a marked agent generates a reward for the punisher of
 211 15 instead of a loss of 20 (note that punishment is always net
 212 negative for the collective reward of the group). We consider
 213 the environment in three conditions corresponding to three
 214 different classification schemes. In the *no rules* condition, no
 215 berries are designated as taboo. Agents never become marked

and punishing is never profitable. In the *important rule* condi-
 216 tion the poisonous berry is taboo. In the *important+silly rule*
 217 condition, both the poisonous berry and another, harmless,
 218 berry are taboo. We manipulate the classification scheme to
 219 assess its causal effect on learning dynamics. We hypothesize
 220 that overall returns are improved by adding the important
 221 rule, and are further improved by adding the silly rule.
 222

223 An agent in our environment has no prior knowledge of
 224 game rules or states. It has no model for the classification
 225 scheme or the potential rewards for appropriately-directed use
 226 of the punishing beam. The agent has to learn how its actions,
 227 its observations (raw pixels), and the rewards it receives relate
 228 to each other entirely from scratch. It can do this by learning
 229 representations that allow it to generalize between similar
 230 situations. Given the enormous game-space and the fact that
 231 all agents have incomplete information about the workings
 232 of the environment, classical game-theoretical analysis is not
 233 tractable. Further, this approach allows us to analyze learning
 234 dynamics, not just equilibria (31).

235 In this framework, behavior is driven by individuals learn-
 236 ing to maximize the expected value of all future rewards they
 237 will obtain from their environment (e.g. by collecting berries,
 238 avoiding and delivering punishment). This learning over time
 239 is accomplished by incremental adjustment of neural network
 240 weights (38). It generates distributed neural representations
 241 that produce reward-maximizing behavior in response to vi-
 242 sual input of the current situation. Agents learn continuously
 243 while being exposed to episode after episode, inhabiting the
 244 same environment with a population of other agents who are
 245 themselves learning simultaneously. In order to do this effec-
 246 tively, agents need to correctly assign credit to current stimuli
 247 and actions based on subsequent rewards they receive. This
 248 creates a rich dynamic in which every part of a behavior has

249 to be learned, and strategic decisions have to be *implemented*
250 via a behavioral policy. Both the cognitive challenge of correct
251 credit assignment (determining which actions contribute to
252 rewards over time), as well as figuring out how to perform
253 complex action sequences are difficult. The dynamics of how
254 norms are learned and implemented are endogenous to the
255 multi-agent learning model. This leads to a number of impor-
256 tant differences from more abstracted simulations like matrix
257 games. We argue that, by focusing on learning, this com-
258 putational model may be particularly appropriate to model
259 anthropological phenomena like the emergence and impor-
260 tance of social norms. In particular the model creates rich
261 learning dynamics for individual agents as well as groups that
262 could not otherwise be approached:

- 263 1. **Complex action sequences** Punishing other agents’
264 behavior, observing a rule violation or complying with
265 a rule are complex sequences of atomic actions that can
266 look different each time they are performed or observed.
- 267 2. **Skills build on each other** As agents have to learn to
268 implement complex behaviors, we can expect a temporal
269 dependency and sequentiality among these behaviors. For
270 example, for agents to learn to avoid a taboo, agents will
271 first need to learn how to effectively apply punishing, in
272 order to motivate rule compliance.
- 273 3. **Opportunity cost** As agents are driven by maximizing
274 total reward, whether or not an agent engages in social
275 punishing depends on the opportunity cost of the action
276 sequence, the agent’s skill in implementing it, and the re-
277 ward gained by punishing the other agent’s transgression.
278 This means there is an intrinsic economy to behavior that
279 is bounded by what agents have learned.
- 280 4. **Generalization** Since the social dynamics are learned in
281 neural networks from scratch they afford the opportunity
282 for, or even necessitate, a degree of generalization. In
283 particular, as punishment is identical for the consequences
284 of transgressing against an important or silly rule, there is
285 an opportunity for generalization of enforcement behavior
286 learned from both rules.
- 287 5. **Endogenous errors** As social punishing of silly or im-
288 portant rules is implemented in the same way, a confusion
289 between the two can arise. Similarly, punishing might be
290 misdirected at agents that did not break a social taboo.
291 These costly false-positive incidents provide an intrin-
292 sic counterweight to the development of an indiscrimi-
293 nate social punishing dynamic. Importantly, unlike other
294 frameworks, multi-agent RL does not require us to model
295 mistakes in behavior as random noise (37, 39). Instead,
296 mistakes in multi-agent RL are emergent from the learn-
297 ing dynamics and the inherent difficulty of implementing
298 an effective behavior policy.

299 Results

300 As displayed in Fig. 4, we examine group-level metrics about
301 agent-populations over the trajectory of learning. We plot the
302 average trajectory per condition. As visible in Fig. 4A, the first
303 thing agent populations learn is to reduce the frequency with
304 which unmarked players are punished. Punishing unmarked
305 players is costly to both the punished and the punishing

agent, so it is unsurprising that this behavior does not persist
306 long once actions become less random. As can be seen in
307 Fig. 4F, this rapid initial learning increases the collective
308 return (the sum of rewards gained by all agents). Note that
309 the suppression of misdirected punishing happens fastest in the
310 *no rules* condition. This is unsurprising, as in this condition
311 there is no direct incentive to punish any other players at all,
312 because there are no taboos that lead to marked players.
313

The second important learning dynamic is that the number
314 of times marked players get successfully punished initially
315 strongly increases before it decreases (Fig. 4B). We interpret
316 the increase as an improvement in the agents’ skill at enforcing
317 the social norm, i.e. being increasingly skilled at effectively
318 punishing marked agents. As displayed in Fig. 4C, the amount
319 of time agents spend marked is steadily declining. However,
320 taken by itself, this metric does not differentiate between
321 whether this decline is driven by agents becoming better at
322 avoiding rule violation, or whether agents get better at pun-
323 ishing rule breakers and thereby removing their mark. As
324 can be seen in Fig. 4E, the decline of successful punishments
325 coincides with a decline in the number of taboo berries eaten.
326 This shows that there is a sequence in the learned behaviors,
327 as first the social punishing system needs to be successfully
328 implemented before it is possible for agents to learn that they
329 should avoid breaking the social norm.
330

In these two measures (successful punishments and taboo
331 berries eaten) we see the role of the silly rule (one additional
332 taboo berry) most clearly. Early in learning, it is unsurprising
333 that doubling the number of taboo berries leads to a higher
334 number of taboo berries eaten and subsequent punishing. But
335 once these quantities start to decline, they decline more rapidly
336 in the condition with two taboos instead of one and in fact
337 reach a lower level. So, it appears that increased exposure
338 to taboo berries and punishing early leads to more robust
339 learning. This is evident in later stages of learning where
340 agents eat fewer taboo berries in the condition in which there
341 are twice as many.
342

As can be seen in Fig. 4D, in terms of avoiding getting
343 poisoned, having two taboos instead of one consistently leads
344 to better results. Additionally, we can see that the credit
345 assignment problem of avoiding the poisonous berry without
346 the help of a social punishing mechanism is prohibitively hard
347 (the *no rule* condition shows no decrease). However, avoiding
348 poisonous berries is not in itself enough to increase overall
349 returns (Fig. 4F). In this environment, overall returns are
350 positively affected by gathering berries and negatively affected
351 by eating poisonous berries and agents punishing each other
352 (every punishment event creates a net loss for the group).
353 Therefore, in order for the additional silly rule to yield an
354 overall benefit in collective return, it needs to not only help
355 to avoid poisonous berries, but it also needs to not add too
356 much dead-weight loss by either hampering consumption of
357 healthy berries or increasing punishment events. As shown in
358 Fig. 4F, there is an overall benefit on collective return in the
359 intermediate learning stages. We test the difference in collec-
360 tive returns between the *important rule* and *important+silly*
361 *rule* condition in 10 separate timebins. There is a significant
362 benefit of the arbitrary rule condition in the 3rd, 4th and
363 5th timebin (3rd: $t(28)=3.94$, $p=0.0005$, 4th: $t(28)=3.26$,
364 $p=0.003$, 6th: $t(28)=2.43$, $p=0.022$. The 3rd and 4th timebin
365 remain significant after Bonferonni-correction for 10 multiple
366

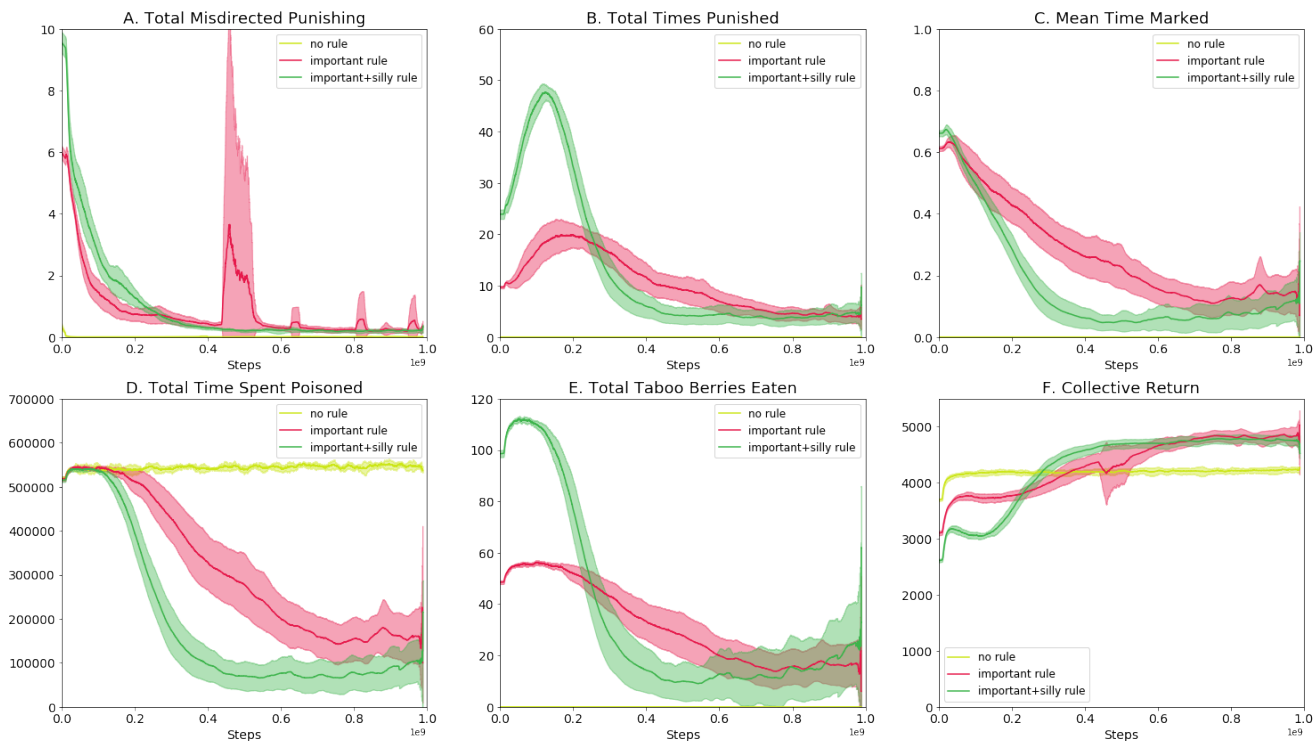


Fig. 4. Learning dynamics: We are examining group-level metrics about agent-populations (y-axis) over the trajectory of learning (x-axis in timesteps). We plot the average trajectory per condition (with 99% confidence interval). A. Number of times unmarked agents are punished (agents that have not broken a taboo). B. Number of times marked agents are punished (agents that have broken a taboo). C. Time spent marked after breaking a taboo. D. Time agents spent poisoned (timesteps after eating the first poisoned berry). E. The number of “taboo” berries eaten (poisonous and non-poisonous combined, if available in the condition). F. Total sum of reward gained by group (including costs of punishing). In total, we observe a benefit of the *important+silly rule* condition in the intermediate stages of learning, driven by an increased ability to avoid poisonous berries. We also see a temporal order to learned behaviors, e.g. an increase in social punishment that then declines together with a decrease in number of taboo berries eaten.

367 comparisons).

368 The results in Fig. 4 suggest that more frequent punishment
 369 early in learning is associated with less time spent poisoned
 370 in the middle stages of learning. We directly test this hypothesis
 371 by exploiting the variance across different training runs (i.e.
 372 separate populations, Fig. 5). In both conditions we find that
 373 high rates of punishment in the early stages of learning (mean
 374 over time, timesteps 0 to 2e8) are related to low amounts of
 375 time spent poisoned in subsequent stages (timesteps 2e8 to
 376 4e8) (*important rule*: $r = -0.79$, $p = 0.0004$; *important+silly*
 377 *rule*: $r = -0.5$, $p = 0.057$, $n = 15$). Note that the correlation
 378 is lower in the *important+silly rule*, but that the magnitudes
 379 of both measures differ strongly. It is possible the correlation
 380 is less pronounced because adding the silly rule increases the
 381 magnitude and restricts range of the rate of punishment (on
 382 the y-axis).

383 **Probing what agents have learned.** Large-scale observational
 384 longitudinal studies with multiple actors face the problem
 385 that all actions taken are entangled and interdependent be-
 386 cause agents react to other agents. Studying the effects of
 387 multiple agents’ interactions over time allows us to investigate
 388 the effects of social norm enforcement on the population at
 389 large, but does not enable conclusions about what specific
 390 mechanisms cause an individual agent’s behaviour. For humans,
 391 psychology experiments address this issue by isolating
 392 specific mechanisms and testing these in controlled conditions,
 393 such as testing reactions to particular stimuli in laboratory

394 experiments. Our simulation allows us to follow this logic and
 395 confront our artificial agents with tightly controlled experimen-
 396 tal environments inspired by lab-testing to directly probe what
 397 the agents have learned. As shown in Fig. 6A, we implement
 398 these quasi-lab experiments by extracting agents at different
 399 points in training and recording their actions when placed in a
 400 simple empty environment with no other agents, and only one
 401 stimulus to interact with. Critically, the agent is not learning
 402 in this environment. Running this experiment with multiple
 403 copies of the same agent allows us to run multiple trials to
 404 probe an agents’ response to a particular game object in isola-
 405 tion. This tests what the agent has learned at different stages
 406 of training. Even though these tests constitute environments
 407 that the agent has not seen during training, the behavioral
 408 results align with what the agent is expected to learn in its
 409 training environment. This is particularly interesting because
 410 it requires a degree of generalization from the agents (‘zero
 411 shot’, as the agents do not learn during the probe). Their
 412 successful transfer of behaviors learned in a large complex
 413 environments to an empty testing environment indicates that
 414 they learned robust behavioral responses to game objects.

415 Fig. 6 B, C & D displays how many timesteps it takes
 416 agents to approach different berries when confronted with
 417 the berries in isolation. These approach-behaviors vary over
 418 the course of training and by the conditions the agents have
 419 learned in. As expected (cf. Fig. 4D), agents learn to avoid
 420 the poisonous berry (pink) in *important rule* and *important+silly*
 421 *rule* but not in *no rule*. Additionally, agents learn to avoid

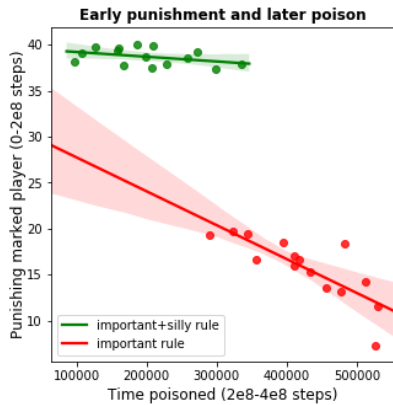


Fig. 5. Higher rates of early punishment is related to less time spent poisoned later in training. Each marker is an independent population of agents. On the y-axis we plot how often players are punished early in training (between timesteps 0 and 2e8). On the x-axis we plot the amount of time players spend poisoned subsequently (mean of the timesteps 2e8 to 4e8). The results are consistent with the interpretation that a high peak in punishment early in training is followed by more avoidance of the poisoned berry later.

the non-poisonous taboo berry (green) in *important+silly rule*. Fig. 6E overlays the poisonous berry lines from panel B and D and illustrates that agents learn to avoid the poisonous berry more in *important+silly rule* than *important rule*. Similarly, Fig. 6F illustrates that agents punish marked players more during learning in *important+silly rule* than *important rule*.

Again, we set out to test the hypothesis that learning about punishment early in training is associated with subsequent avoidance of the poisoned berry. Fig. 6G mirrors the results Fig. 5, demonstrating that the single-player probes are consistent with the behavior observed in the multi-agent setting. We correlate the degree to which a probed agent punishes the marked player during the early stages of training (mean of timewindow 0 to 2e8 steps, marked in panel F) with that agent’s subsequent tendency to consume the poisonous berry (mean of timewindow 2e8 to 4e8 steps, marked in panel E). In both conditions we find a negative relationship (*important rule*: $r = -0.86$, $p < 0.0001$; *important+silly rule*: $r = -0.46$, $p = 0.085$, $n = 15$). Again, note that the absolute magnitude of the values differs across conditions; all datapoints in *important+silly rule* are restricted to relatively high punishment and low rates of approaching the poisonous berry. In sum, these results support the conclusions drawn from the full multi-agent simulation: the additional taboo leads to more frequent punishing events earlier during training, which in turn supports agents’ learning to avoid the poisonous berry. Crucially, these results were obtained in a controlled experimental setting that directly probed what agents have learned by observing their reactions to single objects. These results suggest a sequential, social acquisition of skills, explaining why silly rules help agents learn and behave according to meaningful rules.

Discussion

This work contributes a functional account of why human normative systems contain so many silly and arbitrary rules that is grounded in the mechanics of learning within a single group. The presence of silly rules creates the potential for a larger number of norm violations. From the perspective of

an agent learning the skills necessary to effectively enforce their society’s norms the additional violations constitute additional opportunity for practice, and thus promote a faster rate of improvement in their command of the mechanics of third-party punishment. On the compliance side, the rate at which individuals may learn by trial-and-error to avoid violating taboos depends on the enforcement environment they inhabit. When their groupmates implement highly effective third-party enforcement strategies then exploratory taboo violations are punished both swiftly and surely. Since both speed and certainty of reward (or punishment) are factors known to improve trial and error learning (40, 41), highly “effective” compliance policies (i.e. policies that avoid violating taboos) can be learned rapidly under these conditions. On the other hand, when third-party enforcement is ineffective, then exploratory taboo violations frequently go unpunished or their punishment comes only after a substantial delay. Such conditions are known to make trial and error learning very difficult and slow. Enforcement and compliance are asymmetric in the sense that the former is a skill that may be applied without modification to any norm since many of the sub-behaviors involved in third-party punishment are directed toward the violator (e.g. chasing them), not toward the event of the violation itself. Thus they are “transferable skills”, generically applicable to any norm. Compliance, on the other hand, requires learning to recognize for oneself what would constitute a violation. Now consider also that every society contains a certain number of deeply important rules for which ensuring compliance is of paramount importance. The interpretation of our key result is that the functional role of silly rules in human normative systems may (in part) be to help train a society’s ability to comply with important rules. Adding silly rules into a normative system that already contains deeply important rules can be expected to improve the learning of enforcement for all rules, thereby improving the learning of compliance for all rules, including the rules that truly matter.

While this account is consistent with previous findings on the potential benefits of silly rules (8), the present study demonstrates a novel, mechanistic benefit of silly rules where silly rules improve the scale of enforcement practice, causing a concomitant improvement in the learnability of compliance. We interpret this result as indicating that silly rules enrich the information environment for agents that face a learning challenge. In Hadfield-Menell et al. (8) agents faced a challenge of estimating the likelihood that there are enough agents willing to punish rule violations in a group. In our paper, agents face the more fundamental challenge of learning the relationship between the visual information and actions available to them and the consequences of the two in terms of reward. The silly rule enriches this learning environment with more opportunities to learn about the relationship between punishment behavior and the associated reward for the punisher, as well as the negative consequences for the punished. This account is also independent of, but not necessarily inconsistent with, existing explanations centered around in-group/out-group classification and group cohesion (42). In the real world, adding important rules may be difficult, as they require causal insights into how to avoid undesirable outcomes. Silly rules can be created as needed and, if they are not too costly, the normative order may benefit from the practice that violations against silly rules provide.

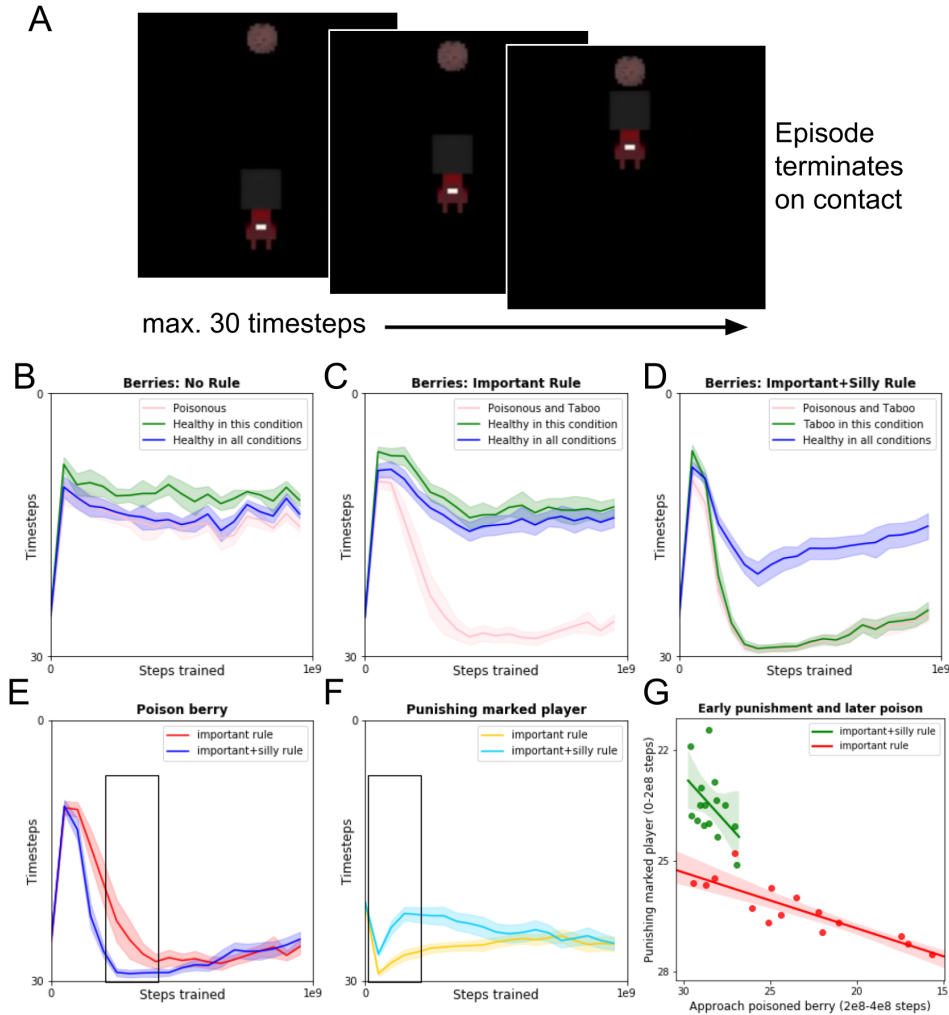


Fig. 6. Single target probes or ‘zero shot generalization’. A. Depiction of probe. An agent is placed in an empty room with just one other object (berry or agent) and we measure how many timesteps it takes to eat the berry or zap the player. B, C & D. Berry types across the 3 different conditions. Agents learn to avoid berries that are taboo. Lines depict the mean across populations of how quickly the agent interacts with the object (y-axis) over learning (x-axis). Error bars represent SEM over different independent populations. E, F. Difference between ‘important rule’ and ‘silly rule’ for approaching the poisoned berry (same as in C & D) and punishing the marked player. Agents are faster to learn to avoid the poisoned berry and punish taboos in the *important+silly rule* condition. G. Early punishing (mean 0 to 2e8 steps) of the marked player is associated with reduced consumption of the poisoned berry (mean 2e8 to 4e8 steps) later in training.

520 While the arbitrary taboo provided a consistent benefit in
 521 avoiding poisonous berries, it is worth noting that the benefit
 522 of the arbitrary rule on the overall prosperity of the group
 523 was only present in the intermediate stages of learning. This
 524 could be associated with the dead-weight cost of maintaining
 525 a social norm that does not serve a direct material function, or
 526 imprecise strategies to avoid poison (i.e. moving more slowly
 527 in general) (43). These costs suggests a strong counterweight
 528 to the usefulness of silly rules in the real world.

529 A clear limitation of this work is that we have not shown
 530 the emergence of the social norms themselves. We supplied in
 531 the environment the causal relationship between an action—
 532 eating a particular berry—and the trigger for social punishing:
 533 becoming marked in the view of other agents and generating
 534 a reward for an agent who successfully aimed the punishing
 535 beam at the transgressor. The next steps in this line of work
 536 are therefore to study the emergence of particular patterns

of marking—norms—and the capacity for norms to change
 in response to changes in the environment or other sources
 or variation including natural drift. We hypothesize that
 learning how to follow and maintain social norms can assist
 agents in adapting to variation in the environment. This social
 technology of benefiting from norms is closely related to the
 cultural niche (11) inhabited by humans, and to humanity’s
 intelligence and success. Understanding how this technology
 emerges in multi-agent settings may play a critical role in
 understanding the emergence of human-level intelligence.

Materials and Methods

Multi-Agent Reinforcement Learning. We consider multi-agent
 reinforcement learning in partially-observable general-sum Markov
 games (44, 45). In each game state, agents take actions based on

552 a partial observation of the state space and receive an individual
 553 reward. Agents must learn through experience an appropriate be-
 554 havior policy while interacting with one another. We formalize
 555 this as follows: an N -player partially observable Markov game
 556 \mathcal{M} defined on a finite set of states \mathcal{S} . The observation function
 557 $\mathcal{O} : \mathcal{S} \times \{1, \dots, N\} \rightarrow \mathbb{R}^d$, specifies each player’s d -dimensional view
 558 on the state space.

559 In each state, each player i is allowed to take an action from its
 560 own set \mathcal{A}^i .

561 Following their joint action $(a^1, \dots, a^N) \in \mathcal{A}^1 \times \dots \times \mathcal{A}^N$, the state
 562 changes obeys the stochastic transition function

563 $\mathcal{T} : \mathcal{S} \times \mathcal{A}^1 \times \dots \times \mathcal{A}^N \rightarrow \Delta(\mathcal{S})$, where $\Delta(\mathcal{S})$ denotes the set of
 564 discrete probability distributions over \mathcal{S} , and every player receives
 565 an individual reward defined as

566 $r^i : \mathcal{S} \times \mathcal{A}^1 \times \dots \times \mathcal{A}^N \rightarrow \mathbb{R}$ for player i . Finally, let
 567 $o^i = \{\mathcal{O}(s, i)\}_{s \in \mathcal{S}}$ be the observation space of player i .

568 Each agent learns, independently through its own experience
 569 of the environment, a behavior policy $\pi^i : \mathcal{O}^i \rightarrow \Delta(\mathcal{A}^i)$ (written
 570 $\pi(a^i | o^i)$), based on its own observation $o^i = \mathcal{O}(s, i)$ and extrinsic
 571 reward $r^i(s, a^1, \dots, a^N)$. Each agent’s goal is to maximize a long
 572 term γ -discounted payoff defined as follows:

$$573 \quad V_{\pi}^i(s_0) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r^i(s_t, \vec{a}_t) \mid \vec{a}_t \sim \pi_t, s_{t+1} \sim \mathcal{T}(s_t, \vec{a}_t) \right]. \quad [1]$$

574 **Experiment and conditions.** We consider a foraging task imple-
 575 mented as a partially observable Markov game on a 2D grid (see
 576 Fig. 2). Agents gain reward by collecting berries that stochasti-
 577 cally respawn. The respawn probabilities are high, so there is little
 578 competition for resources. Moving onto the coordinates of a berry,
 579 agents earn a reward of 4 points. Each berry type is consistently
 580 mapped to a color (24 different types). One berry type is “poi-
 581 sonous”. There is no other signal of which berry type is poisonous
 582 that is observable to an agent at the time of consumption, except
 583 the color that remains consistent for all episodes. If collected by a
 584 player, this player is “poisoned” after a delay of a fixed number of
 585 timesteps (100 timesteps). Poisoning reduces a player’s ability to
 586 absorb nutrition: after poisoning sets in each subsequent berry the
 587 player collects yields a reward of 1 instead of 4. Besides moving,
 588 agents have in their behavioral repertoire the ability to apply a
 589 “punishing beam”. If successfully targeted at another player, the
 590 user of the beam loses a reward of 20 (the cost of punishing, in
 591 addition to the opportunity cost of time spent aiming and firing the
 592 beam instead of collecting berries) and the punished player loses a
 593 reward of 35*.

594 Each instance of the training regime is initialized in one of
 595 3 different conditions. This is a between-subjects design: each
 596 agent population only experiences one of these 3 conditions. The
 597 conditions differ in the content of the classification scheme that
 598 marks agents if they have broken a taboo. We consider three
 599 conditions: no berry is taboo (*no rule*), the poisonous berry is
 600 taboo (*important rule*), the poisonous berry and one harmless berry
 601 are taboo (*important+silly rule*).

602 In *no rule* there are no additional mechanics to the game beyond
 603 what is described above. Agents have to learn which berry is
 604 poisonous without any additional information.

605 In *important rule* we introduce a group rule against eating the
 606 poisonous berry type. In this condition, a player that eats a poison
 607 berry is “marked”: from the perspective of other agents in the
 608 environment, the marked player changes color. This color change
 609 is not visible to the marked player. This color change implements
 610 the idea that other agents evaluate the consumption behavior of
 611 the agent that has chosen to eat a “taboo” food. This marking
 612 then interacts with the punishing capacity of other agents. If
 613 a marked player is successfully targeted by another player with a
 614 punishing beam, the punishing player gets a reward of 35—effectively
 615 transferring reward from the marked player to the punishing player,
 616 for a net payoff to the punishing player of 15 points (note that when
 617 considering the sum of rewards of the whole group, a successful
 618 punishment results net-loss for the group of 20 points because of

619 the cost of using the punishment beam). Aiming punishment at
 620 a non-marked player is costly to both as in the *no rule* condition.
 621 Once punished, the marking disappears.

622 In *important+silly rule*, we augment the important rule with an
 623 additional silly rule, or arbitrary taboo. Players become marked not
 624 only if they consume the poisonous berry but also if they consume
 625 another designated, but harmless, berry. As in the *important rule*
 626 condition, successful punishing of an agent that has violated the silly
 627 rule by consuming the designated harmless berry earns the punishing
 628 agent a net of 15 points and costs the transgressing agent 35 points.
 629 Thus, from the perspective of the agents, the “important” and “silly”
 630 rules are isomorphic if they have not integrated knowledge of the
 631 actual poisoning dynamic.

632 Note that in these settings classification scheme is implemented
 633 by the environment. We have not modeled the emergence of the
 634 rules in themselves. Agents are incentivised to learn policies that
 635 implement the behaviors of collecting berries, delivering third-party
 636 punishment, and avoiding taboo berries that create a risk of pun-
 637 ishment.

638 **Agent architecture and training method.** Each instance of the train-
 639 ing regime contained a population of 12 learners. The environment
 640 is a gridworld of size 33×12 pixels and agents observe a 15×15 pix-
 641 els RGB window, centered on their current location (note that the
 642 depictions in this paper are higher resolution for display purposes).
 643 On each episode, a subset of learners was drawn without replace-
 644 ment to play in the current episode (8 players in each episode). Each
 645 episode lasted for 1000 steps. For each timestep s , each learner i in
 646 the population produced a policy π^i and an estimate of the value
 647 $V_{\pi}^i(s)$ with a neural network, implemented on a GPU. This neural
 648 network was trained by receiving importance-weighted policy up-
 649 dates (46) sampled from a queue of trajectories. These trajectories
 650 were created by 64 simultaneous environments on CPUs that play
 651 the game (with 8 players, which used policies sampled uniformly
 652 from the population of learners without replacement). The learners
 653 received truncated sequences of 100 steps of trajectories in batches
 654 of 16.

655 The neural network’s architecture consisted of a visual encoder
 656 (2D-convolutional neural net with 6 channels, with kernel size and
 657 stride size 1) followed by a 2-layer fully connected MLP with 64
 658 RELU-neurons in each layer, an LSTM (128 units) and finally linear
 659 policy and value heads, outputting the value of the current state and
 660 a probability over actions to be chosen. We used a discount-factor
 661 of 0.99, the learning rate was 0.0004, and the weight of entropy
 662 regularisation of the policy logits was 0.003. We used the RMS-prop
 663 optimiser (learning rate=0.0004, epsilon=1e-5, momentum=0.0,
 664 decay=0.99). The agent also minimized a CPC loss (47) in the
 665 manner of an auxiliary objective (48).

666 **Statistical analysis of observational data.** In order to assess the dif-
 667 ference between conditions, we divide the learning timecourse into
 668 10 bins and average the collective returns for each instance of agent
 669 populations in each bin. We use a t-test to compare the *important*
 670 *rule* and *important+silly rule* conditions in each bin. We correct
 671 the results with a Bonferonni-correction for 10 multiple comparisons
 672 (10 timebins).

673 For the *important rule* and *important+silly rule* conditions we
 674 extracted the mean values for each population of early punishment
 675 (mean of the timesteps 0 to 2e8) and subsequent (mean of the
 676 timesteps 2e8 to 4e8) time spent poisoned. These two measures
 677 were then correlated within each condition.

678 Note that all statistics are done with the datapoints correspond-
 679 ing to entire populations that each contain 12 agents. This is done
 680 because only the data of the entire populations is independent of
 681 each other (the agents within one population affect each other,
 682 therefore do not produce independent data).

683 **Probe methods.** For each agent in each population, the agent’s
 684 unique neural networks were loaded from 20 evenly spaced time-
 685 points spanning the training run. The agent was then placed in a
 686 small empty black environment that contained only one sprite placed
 687 in front of the agent (the sprite of a berry or agent). Each episodes
 688 episode terminates when the agent interacts with the sprite, or after
 689 30 timesteps (timeout). Valid interactions with sprites are “eating”
 690 (upon contact) when the sprite is a berry, and “zapping” with the

* Video of example episode: <https://youtu.be/Xn2eTSX-4GU>. Consumption of taboo berry and sub-
 sequent punishment at 23-25 seconds. Note that agents see a lower resolution version of the
 environment in which each entity is represented by a single pixel.

691 punishment beam when the sprite is an agent. The duration of an
 692 episode is our metric for measuring the agent’s tendency to interact
 693 with the sprite, akin to a “revealed preference” for interacting with a
 694 game object. Shorter episode duration indicates a higher preference
 695 of the agent to interact with the sprite. Note that the agents do
 696 not learn in these episodes. In these probe-episodes, agents are
 697 exposed to the sprite of the pink poisonous berry, a green berry that
 698 is taboo in the *important+silly rule* condition, four berries that are
 699 neither poisonous nor taboo, and the sprite of the “marked” player.
 700 The 20 samples per agent from different timepoints during training
 701 are probed individually with each sprite. Each probe is repeated
 702 20 times and the duration of all episodes is averaged The results
 703 for each timepoint are then averaged across all 12 agents in the
 704 population, resulting in 20 datapoints of each population’s probe
 705 performance over the course of training (15 each for *important rule*
 706 and *important+silly rule* and 5 for *no rule*).

707 **Statistical analysis.** Mirroring the observational data, we extracted
 708 the mean values for each population of early punishment (mean of
 709 the timesteps 0 to 2e8) and subsequent (mean of the timesteps 2e8
 710 to 4e8) approach of the poisoned berry for the *important rule* and
 711 *important+silly rule* conditions. These two measures were then
 712 correlated within each condition.

713 **ACKNOWLEDGMENTS.** The authors would like to thank Jane
 714 X. Wang for helpful discussion.

715 1. PJ Richerson, R Boyd, *Not by genes alone: How culture transformed human evolution.* (Uni-
 716 versity of Chicago press), (2008).
 717 2. K Riedl, K Jensen, J Call, M Tomasello, No third-party punishment in chimpanzees. *Proc.*
 718 *Natl. Acad. Sci. United States Am.* **109**, 14824–14829 (2012).
 719 3. M Tomasello, A Vaish, Origins of human cooperation and morality. *Annu. Rev. Psychol.* **64**,
 720 231–255 (2013).
 721 4. JW Buckholtz, R Marois, The roots of modern justice: cognitive and neural foundations of
 722 social norms and their enforcement. *Nat. neuroscience* **15**, 655–661 (2012).
 723 5. DM Fessler, CD Navarrete, Meat is good to taboo. *J. Cogn. Cult.* **3**, 1–40 (2003).
 724 6. J Henrich, N Henrich, The evolution of cultural adaptations: Fijian food taboos protect against
 725 dangerous marine toxins. *Proc. Royal Soc. B: Biol. Sci.* **277**, 3715–3724 (2010).
 726 7. MC Corballis, Laterality and myth. *Am. Psychol.* **35**, 284 (1980).
 727 8. D Hadfield-Menell, M Andrus, G Hadfield, Legible normativity for ai alignment: The value
 728 of silly rules in *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*,
 729 (ACM), pp. 115–121 (2019).
 730 9. G Brennan, L Eriksson, RE Goodin, N Southwood, *Explaining norms.* (Oxford University
 731 Press), (2013).
 732 10. R McElreath, R Boyd, P Richerson, Shared norms and the evolution of ethnic markers. *Curr.*
 733 *anthropology* **44**, 122–130 (2003).
 734 11. R Boyd, PJ Richerson, J Henrich, The cultural niche: Why social learning is essential for
 735 human adaptation. *PNAS* **108**, 10918–10925 (2011).
 736 12. B Kenward, Over-imitating preschoolers believe unnecessary actions are normative and en-
 737 force their performance by a third party. *J. experimental child psychology* **112**, 195–207
 738 (2012).
 739 13. GK Hadfield, BR Weingast, Microfoundations of the rule of law. *Annu. Rev. Polit. Sci.* **17**,
 740 21–42 (2014).
 741 14. M Derex, JF Bonnefon, R Boyd, A Mesoudi, Causal understanding is not necessary for the
 742 improvement of culturally evolving technology. *Nat. human behaviour* **3**, 446 (2019).
 743 15. R Bond, PB Smith, Culture and conformity: A meta-analysis of studies using asch’s (1952b,
 744 1956) line judgment task. *Psychol. bulletin* **119**, 111 (1996).
 745 16. J Henrich, R Boyd, The evolution of conformist transmission and the emergence of between-
 746 group differences. *Evol. human behavior* **19**, 215–241 (1998).
 747 17. W Hoppitt, KN Laland, *Social learning: an introduction to mechanisms, methods, and models.*
 748 (Princeton University Press), (2013).
 749 18. C Tennie, J Call, M Tomasello, Ratcheting up the ratchet: on the evolution of cumulative
 750 culture. *Philos. Transactions The Royal Soc. B* **364**, 2405–2415 (2009).
 751 19. GK Hadfield, BR Weingast, Law without the state: legal attributes and the coordination of
 752 decentralized collective punishment. *J. Law Court.* **1**, 3–34 (2013).
 753 20. P Richerson, R Boyd, The evolution of human ultra-sociality in *Ideology, Warfare, and Indoc-*
 754 *trinability*, eds. I Eibl-Eibesfeldt, F Salter. (Berghen Books, Oxford), pp. 71–95 (1998).
 755 21. J Henrich, M Muthukrishna, The origins and psychology of human cooperation. *Annu. Rev.*
 756 *Psychol.* **72**, 207–240 (2021).
 757 22. JZ Leibo, E Hughes, M Lanctot, T Graepel, Autocurricula and the emergence of innova-
 758 tion from social interaction: A manifesto for multi-agent intelligence research. *arXiv preprint*
 759 *arXiv:1903.00742* (2019).
 760 23. R Boyd, H Gintis, S Bowles, PJ Richerson, The evolution of altruistic punishment. *Proc. Natl.*
 761 *Acad. Sci.* **100**, 3531–3535 (2003).
 762 24. LC Brooks, W Iba, S Sen, Modeling the emergence and convergence of norms in *Twenty-*
 763 *Second International Joint Conference on Artificial Intelligence.* (2011).
 764 25. S Mahmoud, S Miles, M Luck, Cooperation emergence under resource-constrained peer
 765 punishment in *Proceedings of the 2016 International Conference on Autonomous Agents*
 766 *& Multiagent Systems.* (International Foundation for Autonomous Agents and Multiagent
 767 Systems), pp. 900–908 (2016).

26. N Ajmeri, H Guo, PK Murukannaiah, MP Singh, Robust norm emergence by revealing and
 768 reasoning about context: Socially intelligent agents for enhancing privacy. in *IJCAI.* pp. 28–34
 769 (2018).
 770 27. M Kleiman-Weiner, MK Ho, JL Austerweil, ML Littman, JB Tenenbaum, Coordinate to coop-
 771 erate or compete: abstract goals and joint intentions in social interaction in *CogSci.* (2016).
 772 28. JZ Leibo, V Zambaldi, M Lanctot, J Marecki, T Graepel, Multi-agent reinforcement learning
 773 in sequential social dilemmas in *Proceedings of the 16th Conference on Autonomous Agents*
 774 *and MultiAgent Systems.* (International Foundation for Autonomous Agents and Multiagent
 775 Systems), pp. 464–473 (2017).
 776 29. A Lerer, A Peysakhovich, Maintaining cooperation in complex social dilemmas using deep
 777 reinforcement learning. *arXiv preprint arXiv:1707.01068* (2017).
 778 30. A Peysakhovich, A Lerer, Consequentialist conditional cooperation in social dilemmas with
 779 imperfect information. *arXiv preprint arXiv:1710.06975* (2017).
 780 31. J Perolat, et al., A multi-agent reinforcement learning model of common-pool resource appro-
 781 priation in *Advances in Neural Information Processing Systems.* pp. 3643–3652 (2017).
 782 32. E Hughes, et al., Inequity aversion improves cooperation in intertemporal social dilemmas in
 783 *Advances in neural information processing systems.* pp. 3326–3336 (2018).
 784 33. J Foerster, et al., Learning with opponent-learning awareness in *Proceedings of the 17th*
 785 *International Conference on Autonomous Agents and MultiAgent Systems.* (International
 786 Foundation for Autonomous Agents and Multiagent Systems), pp. 122–130 (2018).
 787 34. A Peysakhovich, A Lerer, Prosocial learning agents solve generalized stag hunts better than
 788 selfish ones in *Proceedings of the 17th International Conference on Autonomous Agents*
 789 *and MultiAgent Systems.* (International Foundation for Autonomous Agents and Multiagent
 790 Systems), pp. 2043–2044 (2018).
 791 35. JX Wang, et al., Evolving intrinsic motivations for altruistic behavior in *Proceedings of the 18th*
 792 *International Conference on Autonomous Agents and MultiAgent Systems.* (International
 793 Foundation for Autonomous Agents and Multiagent Systems), pp. 683–692 (2019).
 794 36. B Baker, Emergent reciprocity and team formation from randomized uncertain social prefer-
 795 ences. *Adv. neural information processing systems (NeurIPS)* (2020).
 796 37. R Boyd, S Mathew, Arbitration supports reciprocity when there are frequent perception errors.
 797 *Nat. Hum. Behav.*, 1–8 (2021).
 798 38. Y LeCun, Y Bengio, G Hinton, Deep learning. *nature* **521**, 436–444 (2015).
 799 39. RB Myerson, Refinements of the nash equilibrium concept. *Int. journal game theory* **7**, 73–80
 800 (1978).
 801 40. RS Sutton, AG Barto, A temporal-difference model of classical conditioning in *Proceedings*
 802 *of the ninth annual conference of the cognitive science society.* (Seattle, WA), pp. 355–378
 803 (1987).
 804 41. W Dabney, et al., A distributional code for value in dopamine-based reinforcement learning.
 805 *Nature* **577**, 671–675 (2020).
 806 42. VB Meyer-Rochow, Food taboos: their origins and purposes. *J. Ethnobiol. Ethnomedicine* **5**
 807 (2009).
 808 43. E Nisioti, C Moulin-Frier, Grounding artificial intelligence in the origins of human behavior.
 809 *arXiv preprint arXiv:2012.08564* (2020).
 810 44. LS Shapley, Stochastic Games. In *Proc. Natl. Acad. Sci. United States Am.* (1953).
 811 45. ML Littman, Markov games as a framework for multi-agent reinforcement learning in *Pro-*
 812 *ceedings of the 11th International Conference on Machine Learning (ICML).* pp. 157–163
 813 (1994).
 814 46. L Espeholt, et al., Impala: Scalable distributed deep-rl with importance weighted actor-learner
 815 architectures (2018).
 816 47. Avd Oord, Y Li, O Vinyals, Representation learning with contrastive predictive coding. *arXiv*
 817 *preprint arXiv:1807.03748* (2018).
 818 48. M Jaderberg, et al., Reinforcement learning with unsupervised auxiliary tasks. *arXiv preprint*
 819 *arXiv:1611.05397* (2016).
 820