

Uncertainty, bias and the institution bootstrapping problem

Stavros Anagnou^{1,2}, Christoph Salge¹, and Peter R. Lewis²

¹ Adaptive Systems Research Group, University of Hertfordshire, United Kingdom

² Trustworthy AI lab, Ontario Tech University, Canada

`s.anagnou@herts.ac.uk`

Abstract. Institutions play a critical role in enabling communities to manage common-pool resources and avert tragedies of the commons. Prior research suggests institutions emerge when universal participation yields greater collective benefits than non-cooperation. However, a fundamental issue arises: individuals typically perceive participation as advantageous only after an institution is established, creating a paradox—how can institutions form if no one will join before a critical mass exists? We term this conundrum the institution bootstrapping problem and propose that misperception—specifically, agents’ erroneous belief that an institution already exists—could resolve this paradox. By integrating well-documented psychological phenomena—including cognitive biases, probability distortion, and perceptual noise—into a game-theoretic framework, we demonstrate how these factors collectively mitigate the bootstrapping problem. Notably, unbiased perceptual noise (e.g., noise arising from agents’ differing heterogeneous physical or social contexts) drastically reduces the critical mass of cooperators required for institutional emergence. This effect intensifies with greater diversity of perceptions, suggesting that variability among agents’ perceptions facilitates collective action. We explain this counter-intuitive result through asymmetric boundary conditions: proportional underestimation of low-probability sanctions produces distinct outcomes compared to equivalent overestimation. Furthermore, the type of perceptual distortion—proportional versus absolute yields qualitatively different evolutionary pathways. These findings challenge conventional assumptions about rationality in institutional design, highlighting how "noisy" cognition can paradoxically enhance cooperation. Finally, we contextualize these insights within broader discussions of multiagent system design, noise and cooperation, and disobedience and collective action. Our analysis underscores the importance of incorporating human-like cognitive constraints—not just idealized rationality—into models of institutional emergence and resilience.

Keywords: Uncertainty · Bias · Institutions · Bootstrap · Evolutionary Game Theory · Noise · Bounded Rationality

1 Introduction

1.1 The institution bootstrapping problem

Across human societies, it is difficult to find groups without some form of institution [24]. Institutions—defined as rule systems enabling sustainable management of common-pool resources (e.g., grazing lands, fisheries, or even Minecraft servers)—help communities avoid antisocial outcomes by regulating individual self-interest [25, 10]. There are many empirical examples of such institutional rules empowering groups to avoid anti-social outcomes by regulating individual self-interest in order to manage the resource sustainably e.g. how much water an individual may take from a shared irrigation system and how often they must perform maintenance. Critically, institutions are self-reinforcing: both compliance and active enforcement align collective good with individual incentives, as demonstrated in foundational studies [24, 25].

However, institutions require sustained effort to create, maintain, and adapt. Without this, groups risk reverting to default interactions where cooperation collapses. Ostrom’s fieldwork emphasizes that institutions are more likely to endure when rules are both designed and enforced by the same agents whose actions they govern [25]. Institutional roles—such as monitoring compliance or coordinating rule updates—incur costs that must be offset by collective benefits in order for self-interested individuals to join. Recent work formalizes these dynamics [27, 28]. For example, Powers et al [28] use evolutionary game theory (EGT), to derive conditions where institutional participation becomes evolutionarily stable. For instance, monitoring costs must remain low relative to resource contribution costs, and institutional roles (e.g., monitors) must be incentivized through resource-sharing mechanisms [28].

A critical challenge arises even when institutions are theoretically favourable: the institution bootstrapping problem. As Powers et al [28] notes, institutions require a threshold of participants to generate sufficient benefits (e.g., monitoring capacity) to justify individual costs. Without this critical mass, free-riding dominates because early adopters bear disproportionate costs (e.g., monitoring efforts) without guaranteed reciprocity.

So, if we assume all individuals in a group are not yet part of the institution, the benefits do not outweigh the costs of joining. Since there are no monitors to enforce rules on peers and therefore a lesser incentive to contribute as opposed to free-riding. We term this issue the institution bootstrapping problem. We consider several explanations for this and possible approaches for alleviating or completely overcoming the problem, where agents would manage to bootstrap their way towards the institution despite the incentives against it.

Proposed solutions include extrinsic shocks (e.g., external cooperator influx) or cognitive factors like cost misperception [28]. Small groups might circumvent the problem through charismatic leadership or evolved psychological mechanisms that amplify cooperation e.g. trust [18]. Our approach focuses on relaxing perfect rationality assumptions in evolutionary game theory by incorporating perceptual

biases and uncertainty in the form of noise inherent to bounded agents operating in heterogeneous physical/social environments.

1.2 Contributions

Using an evolutionary game theoretic approach we:

1. Explicitly illustrate, for the first time, the bootstrapping problem in a simplex plot.
2. Show that incorporating a coarse bias into how agents perceive the cost of freeriding, can either decrease or increase the number of cooperators and monitors needed to establish the institution. With overestimating the risk of punishment leading to greater cooperation. This holds if we generally assume that the agents are subject to a loss aversion bias, but given complicated emerging empirical evidence on how individuals perceive probabilities we need a more solid empirical foundation to motivate this bias.
3. Show that incorporating a more nuanced S-shaped or inverse-S-shaped probability distortions from the literature on human psychophysics, we can similarly show that this effects the threshold of cooperators needed to establish an institution. However, limitations due to the individual context of the experiments and task dependency of the effect complicate interpretation. We, therefore, urge for such psychophysics experiments to take place in a social context as their implications would matter for the psychology of enduring institutions.
4. Show that incorporating noisy perception to capture the bounded heterogeneity among agents (e.g. their differing social circles and positions in the physical world) decreases the threshold of cooperators needed to establish the institution. Interestingly, this is despite the noise being unbiased at the group level.
5. Explain the above counter-intuitive result in terms of an asymmetric boundary condition where proportionally underestimating very small quantities is not the same as overestimating them and show that the type of noisy perception (proportional or absolute), results in different qualitative results.

2 Modelling the institution bootstrapping problem

In this section, we will describe the bootstrapping problem in terms of evolutionary game theory (EGT).

To capture the bootstrapping problem in more concrete terms, we will express it in the form of a game-theoretic model. Investigating such theories in a mathematical framework helps us explicitly define our assumptions and system specifications. It allows us to assess the logic of our ideas and establish whether they are internally coherent enough to serve as a good candidate explanation. This added rigour helps us avoid logical errors or missed details due to the inherent ambiguity of verbal theorizing.

Further, the abstractions in a mathematical model allow us to capture dynamics common across many types of institutions and provide a general understanding of how they work without getting lost in details. This gives us a solid theoretical foundation to subsequently incorporate the specifics of each situation.

We adapt the Powers model of institutions [28], using it to describe explicit utility functions within an evolutionary game theory (EGT) framework for each action [34]. EGT evaluates the utility of an agent’s actions, which reflects the material or psychological consequences of those actions. These utility functions are then combined with a replicator equation to plot the rate of change in the relative frequencies of strategies, generating a simplex plot that visualizes population dynamics across states (i.e., the number of agents adopting each strategy).

We outline three strategies available to individuals: defector (D), contributor (C), and contributor-monitor (CM). A **defector (D)** does not engage in the institution, undermining it by consuming from the common-pool resource without contributing (C_c) and avoiding institutional roles. Defectors gain benefits from the common pool (B_g , Table 1) but incur a freeriding cost (C_f , Table 1) in the form of punishment from monitors. This cost scales with N_m/N , the fraction of monitors in the population, reflecting the increased likelihood of being caught as the proportion of monitors rises (Table 2). It also scales with p , the number of checks a monitor makes and s , the cost of punishment. This weighted expectation of being punished (C_f , as captured abstractly by a reduction in utility, models the potential consequences of peer punishment e.g. a material fine, the feeling of shame or damage to reputation) [32, 23]. This deters individuals from defecting and maintains social order [23].

The **contributor (C)** participates in the institution, benefiting from it (B_g , Table 2) while paying a contribution cost (C_c , Table 2) determined by the parameter α , which quantifies the individual’s contribution level (Table 2).

The **contributor-monitor (CM)** also participates and contributes but additionally takes on a monitoring role to enforce institutional rules. Like contributors, CM agents benefit (B_g) and pay contribution costs (C_c). Monitoring incurs an additional cost (C_m , Table 2), which includes effort (e.g., time spent checking compliance) or risks (e.g., retaliation from punished defectors). These costs are modelled as $p \cdot \delta$, where p represents the number of monitoring checks performed, and δ is the cost per punishment instance. Monitors receive a benefit (B_m , Table 2) proportional to β (the share of common-pool resources allocated to monitors), α (individual contributions), and N_c/N_m (the contributor-to-monitor ratio), with larger N_m diluting individual shares.

We have described a system where the consequences of an agents’ actions depend on facts in the world e.g. the cost of punishment s and the frequency of other strategies in the population e.g. the amount of monitors determines how much an agent is punished (Tables 1,2,3). This means the utility of each individual depends on what other individuals are doing.

The utilities of each strategy (D,C,CM) determine the strategy’s success and their propagation through the population. In classic EGT this is cast in terms of genetic evolutionary fitness. But we can also interpret this to be cultural

Table 1. Table outlining strategies in terms of their role, real world example and utility payoff in model

Strategy	Equation	Role	Real world examples (local town)
D	$U_D = B_g - C_f$	Agent who is depleting common pool resource but not contributing to it, is punished by monitors	Individual uses public infrastructure but doesn't pay to upkeep it
C	$U_C = B_g - C_c$	Agent contributing to common pool resource but not taking on a monitoring role to punish D	Individual uses public infrastructure and pays to upkeep it
CM	$U_{CM} = B_g - C_c + B_m - B_c$	Agent contributing to common pool resource and taking on a monitoring role to punish D	Individual uses public infrastructure, pays to upkeep it and goes out of their way to punish others who do not

Table 2. Table describing the associated costs and benefits of each action along with an equation defining it.

Cost or Benefit	Meaning	Equation
B_g	Benefit from collective resource of being a member of the group	$B_g = (1 - \beta) \cdot 1/N \cdot \alpha \cdot N_c$
C_c	Cost of contributing	$C_c = \alpha$
B_m	Benefit of monitoring	$B_m = \alpha \cdot \beta \cdot N_c / Nm$
C_m	Cost of monitoring	$C_m = p \cdot \delta$
C_f	Cost of free-riding (note that this is a expected cost conditioned on a probability)	$C_f = (p \cdot N_m / N) \cdot s$

Table 3. Table describing each world parameter that is used to compute the utility of each strategy

Parameter	Meaning	Value (to favour institution)	Value (to not favour institution)
α	Cost of contributing to common pool resource	1	1
β	Fraction of common pool resource given to monitors	0.2	0.2
δ	Cost of punishment to monitor	0.1	0.5
p	Number of checks each monitor makes	5	5
s	Cost of being punished	1	1
N_m	Number of monitors	Dynamic	Dynamic
N_c	Number of contributors (includes CMs)	Dynamic	Dynamic
N	Total number of agents	20	20

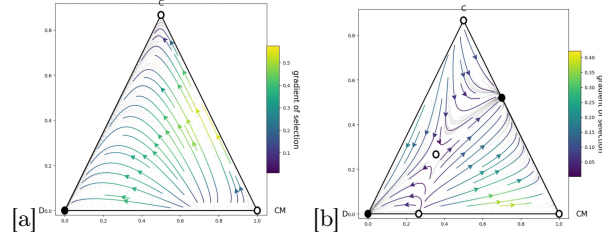


Fig. 1. (a) depicts the simplex plot for 3 strategies where the inequality specified by Powers et al [28] (a mathematical inequality dependent on world parameters that determines whether or not an institution will be formed) is not satisfied and an institution does not form (b) depicts a case where the inequality from Powers [28] where an institution does form, however for areas in strategy space closer to all defectors, the system tends toward all defectors. This means although conditions in terms of parameters are favourable for an institution to form, a critical mass of individuals joining the institution is needed before it will form. Reaching this threshold is the institution bootstrapping problem. The black circles depict unstable attractors, and the black dots signify stable attractors

fitness i.e. strategies that are better off would be more likely to be copied by others through prestige biased social learning [4]. We implement this using the replicator equation (equation 1).

The replicator equation computes the gradient, for a given group configuration x_i . It does this by multiplying the average fitness of a strategy across at the population at a given group configuration $\Pi_i(\mathbf{x})$ vs the average fitness of all strategies for that given group configuration $\bar{\Pi}(\mathbf{x})$ (equation 1).

$$\dot{x}_i = x_i (\Pi_i(\mathbf{x}) - \bar{\Pi}(\mathbf{x})) \quad (1)$$

We can use the replicator equations to plot the direction in strategy space for each group combination in a simplex and show when agents will tend to the pro-institution strategies (C and CM) and then they will not (D). We do this using the library EGTtools [9].

Powers derived conditions under which an institution will endure, expressed through mathematical inequalities [28]. Specifically, these inequalities compare institutional benefits against participation costs: an institution endures when benefits exceed costs.

In Fig. 1, we plot simplexes representing strategy space gradients for systems that do and do not satisfy these inequalities, illustrating the bootstrapping problem. Arrows indicate selection gradients, showing the direction of strategy propagation under evolutionary dynamics. The corners (D, C, CM) correspond to homogeneous populations (all agents adopting one strategy), while the centre represents equal strategy proportions. White dots denote unstable equilibria (system states persisting only without stochastic changes, e.g., imitation errors during strategy updates). Black dots signify evolutionary stable strategies (ESS),

robust to perturbations in group composition and therefore cannot be invaded by another strategy.

In Fig. 1a, the simplex corresponds to parameters predicted to preclude institutional formation per [28]. Here, institutions fail to form in almost all configurations, with no stable attractor for populations of C or CM. The sole stable attractor occurs at full defection (D), indicating an institution will not form in almost any case and if it does it will be fragile against invasion by non-cooperative strategies.

Fig. 1b shows parameters predicted to enable institutional formation per [28]. A large region of the strategy space converges toward an ESS for institutional cooperation (C/CM dominance), where defection yields no advantage. This ESS resists large perturbations in group composition. However, a smaller region persists where insufficient monitors/cooperators render cooperation non-beneficial, resulting in a D-dominated ESS. In this region, isolated strategy shifts (e.g., one agent cooperating) revert to D due to lack of collective incentives.

Thus we have captured the institution bootstrapping problem: despite the institution being beneficial to join, if the system starts in an initial non-cooperative situation (which is an equilibrium), agents need a catalyst to be able to bootstrap the creation of institutions, which can then lead to socially preferable, sustainable outcomes (also shown to be an equilibrium).

3 Incorporating perception into social simulation, an unlikely solution to the bootstrapping problem?

How do we get around the bootstrapping problem? One can assume an influx of cooperators/monitors to the group which then incentivises others to join, or a strong leader type emerges which forces the requisite number of individuals to join to incentivise institution formation [28]. However, an unexplored avenue in addressing this problem may come from incorporating facets of human psychology into social modelling, which would then allow us to try an unorthodox approach: Can agents merely pretending an institution exists, make it a reality?

We will now try to alleviate the bootstrapping problem by questioning the cognitive assumptions of EGT.

As seen in the equation for the expected cost of freeriding C_f (Table 2), it assumes a perfect perception of what would often be hidden variables e.g. the number of checks a monitor makes p or the amount of monitors in the population at any given time.

It is very difficult for bounded agents with noisy perception to have a perfect estimate, in most cases they often resort to heuristics in the form of biases [36, 33, 13]. Emerging work in game theory and social simulation [20, 37], shows that bringing even simple psychological facets in can change the predictions of game theory and the dynamics of agent based models [20].

We focus on bias in the literature of social modelling and evolutionary game theory.

Firstly there is work showing that bias, even though it may be misrepresentative of reality, may be advantageous to individuals. For example, in a frequency-dependent hawk-dove game, where the prevalence of doves and hawks influences strategic payoffs, overestimating the reward of a hawk action can mislead unbiased individuals into perceiving their own actions as less advantageous. As a result, they alter their behaviour, indirectly benefiting biased hawk-strategy players, who now face fewer competitors [21]. Similarly, in resource competition, individuals who are overconfident in their ability to compete are evolutionarily stable across a wide range of environments, particularly under uncertainty [14].

In the above cases, even though biased individuals were better off, bias was detrimental for group utility, but there are also cases where it is beneficial for the group as well. For example, Vogrin et al [37] show that having a bias enhances performance in a signal detection task and suggest that in multiple agent settings, different agents can do a cognitive division of labour by specialising in different signals thus promoting discourse and a broader investigation of problem spaces. Therefore being beneficial to the whole group. Davies et al [6] show that incorporating an adaptive bias in how one’s utility is perceived in a coordination problem can change the attractor landscape of a system, widening the basin of the global optimum and therefore making it more likely for collective systems to arrive at the global optima and therefore benefit all. We hope to apply similar tactics to enlarge the cooperative attractor/reduce the number of individuals needed to set up an institution, in our scenario.

4 Experiments

We modify the expected cost of freeriding C_f in various ways by biasing the value away from the perfect expected probability to model cognitive biases and aspects of an agent’s perceptual uncertainty (noise) and boundedness.

1. Coarse grained bias: bias perception of C_f (expected cost of freeriding) by multiplying by 0.75 or 1.5 (for under or over-estimating)
2. Distorting perception of extreme probabilities: The inverse S and S shaped probability distortion curves are modelled by the Prelec function, given by $\pi(p) = e^{-\zeta(-\ln p)^\lambda}$, where $\pi(p)$ represents the subjective probability, and p is the objective probability within the range ($0 < p \leq 1$). The parameter λ controls the curvature of the function: if $\lambda < 1$, the function follows an inverted-S shape (we used 0.8), meaning small probabilities are overweighted while large probabilities are underweighted; if $\lambda > 1$, (we used 1.2) the function follows an S-shape, where small probabilities are underweighted and large probabilities are overweighted. The parameter ζ adjusts the elevation of the curve, typically set to 1 in standard applications.
3. Noise (proportional), here we sample equally above >1 and <1 when doing the fitness calculations in order to simulate a noisy perception as would occur with agents bounded by individual and social context. $noise \sim \mathcal{U}(a, b)$ which is then multiplied by C_f

4. Noise (absolute): Here we do the same as 3 but we add the noise so it is absolute (we add or subtract in the given range).

4.1 Moran process for simulating finite populations

To model these effects of noisy bias we need a stochastic payoff for the defect strategy, which is not possible in a replicator equation as they assume infinite populations and are deterministic.

Therefore, we will be using a Moran process, which assumes finite populations and therefore allows for stochastic effects.

We use the library EGTtools to implement the Moran process [9]. For smaller population sizes, stochastic effects dominate, necessitating discrete birth-death processes to model behavioural dynamics. This framework introduces the **finite population selection gradient** $G(k/Z)$, defined as the difference between probabilities of incrementing or decrementing a strategy's count, stochasticity intensifies with behavioural "mutations" (e.g., imitation errors).

We model social learning dynamics via a stochastic birth-death process paired with a pairwise comparison rule. At each timestep, a randomly selected individual j (strategy j) revises their strategy by potentially imitating a randomly chosen individual i . Imitation probability p follows the **Fermi function**:

$$p = \frac{1}{1 + e^{\gamma(f_i(k_i) - f_j(k_j))}} \quad (2)$$

where $f_i(k_i)$ and $f_j(k_j)$ denote the fitness of individuals i and j , dependent on their strategy abundances k_i and k_j . Due to finite populations, absolute counts k_i replace frequencies x_i , where $x_i \equiv k_i/Z$.

Here, γ (inverse temperature) modulates selection intensity and imitation accuracy: $\gamma \rightarrow 0$ induces near-random drift; $\gamma \rightarrow \infty$ renders imitation deterministic. A mutation rate μ allows random strategy exploration. Collectively, this adaptive process forms a Markov chain with state transitions governed by strategy fitness and abundance. For specifics, see [9].

5 Results

5.1 Coarse bias

Due to limited cognitive resources, agents usually rely on imperfect but frugal heuristics e.g. loss aversion heuristic [19, 36].

We incorporate such a bias into our model as a proportional bias on the C_f and see if this changes the attractor landscape and helps address the bootstrapping problem.

We see in Fig.2b, where agents overestimate the risk of being punished, that the defector attractor indeed shrinks, making it easier for influxes of cooperators or other extrinsic shocks of the system to be able to meet the critical mass of joiners needed to form an institution. For the sake of completeness we show in

Fig.2c that underestimating the risk leads to a larger defective attractor, and therefore it makes it harder for institutions to form.

It has been shown empirically that humans overestimate unlikely punishments due to loss aversion bias [19, 36], which then would help institutions form. Further, evidence shows that humans overestimate vivid but rare risks e.g. as can be induced by media sensationalism [35]. This could explain why public punishment/moral panics are effective ways to influence people to join institutions. We could then argue this general risk aversion evolved and was then co-opted for institutions. However, more work has been done on bias that complicates this picture of human bias [36, 38]. We will incorporate this work in the next section.

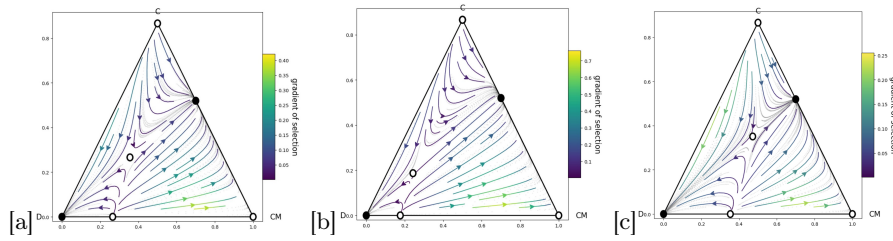


Fig. 2. (a) Control (b) Agents overestimate the expected cost of punishment (c) Agents underestimate the expected cost of punishment

5.2 Bias in the form probability distortion

Empirical work has also shown that humans often distort their perception of probabilities at the extremes in a non-linear manner (i.e. distort very small and very large probabilities) [38]. This results in either an inverted S shaped distortion (overvalue small probabilities and underestimate large ones) or an S shaped distortion (undervalue small probabilities and over estimate large ones). We incorporate S and inverse S distortion from the experimental literature into our model to see what how it effects the likelihood of institution formation.

Under the inverted S distortion, similar to the coarse overestimation bias, we see an increase in the size of the attractor for cooperation, which therefore makes cooperation more likely. When we have an inverted S distortion and a decrease in the size of the cooperative attractor when it is S shaped, which makes cooperation more unlikely (Fig.3). This is interesting as it suggests incorporating these findings from psychophysics into game theory changes collective outcomes.

However, despite its more detailed empirical grounding, the result here is hard to interpret as it hinges on which curve (S or inverse S) humans use when forming institutions. The empirical evidence for this is not clear cut. For example, it seems that use of S or inverse S shaped curves are heterogeneous in population, so not everyone has the same shape. Secondly, it also depends on type of task i.e. motor decision task vs abstract economic decision task [38]. Furthermore, social

contexts can change outcomes e.g. when people are told they are playing a game with people vs a computer their behaviour changes [29] and whether a participant can see if they are playing against a real person or not also affects behaviour [7]. Such experiments on probability distortion have not been attempted in a truly social setting. Therefore understanding the relation between probability distortion and collective outcome is not straightforward.

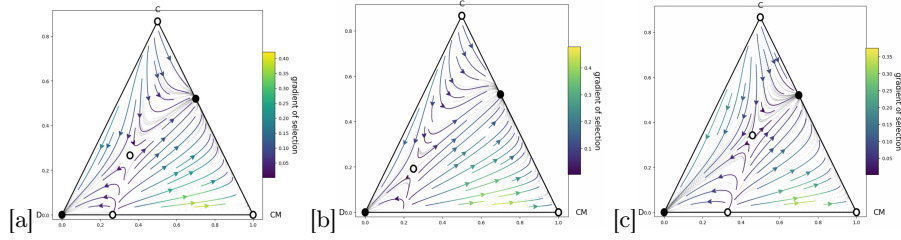


Fig. 3. (a) Control (b) Agents obey an inverted S shape distortion in how they estimate expected punishment (overestimate small probabilities and underestimate large ones) (c) Agents obey an S shape distortion in how they estimate expected punishment (underestimate small probabilities and overestimate large ones)

5.3 Noise (bounded uncertainty)

Given the difficulty of interpretation of empirical results to our institution bootstrapping problem. We now attempt to motivate this in an alternate fashion using the inherent uncertainty brought about because individuals are bounded by their local context. It has been established that individual agents are bounded [33] and often limited to the information around them. Therefore they often anchor on local estimates and recent experiences e.g. people who know someone with cancer have higher estimates of its prevalence [33, 26]. These estimates are subject to noise for a host of different factors e.g. different life experiences, affective make ups of each individual, social networks, social media bubbles e.t.c [35, 15, 26].

In our model, an agent may be very optimistic with respect to being punished based on its life experience, and may have been lucky enough not to have been punished before and be in a social group that either doesn't get punished very often or doesn't like advertising it to others. This would make the agent relatively optimistic with respect to the dangers of punishment (underestimates risk). Alternately, another agent could have been punished quite a few times and be in a social group that gossips about punishment often, leaving it with a rather bleak estimate of how likely it is to get punished (overestimates risk). Because of the uncertainty inherent in being bounded agents have to anchor on local cues which results in a spectrum differing estimates of being punished across the population.

To model the inherent uncertainty of individuals bounded by different contexts and noisy perception, we represent it as a uniform distribution of biases (away from the perfect estimate) across the population. Note that this is not the same as all individuals population having the same bias as in previous sections, here each individual has different bias determined by their local social and physical context. Further, there is no skew in the bias at the population level, all noisy estimates even out to be an accurate unbiased group estimate, as supported by collective intelligence studies [11, 17, 16].

To model this, we need a stochastic payoff for the defect strategy. To do this we sample uniformly for in a number in a range with 40 samples (equal to population) centred at one. Which is then multiplied by C_f so that effectively each agent has their own estimate of the expected cost of freeriding, with equal chances of them under or over valuing the probability of being punished.

We see that, surprisingly, under this noise, cooperative attractors are larger (Fig.4). This doesn't solve the bootstrapping problem completely but means stochastic extrinsic shocks e.g. strong man mutation or an influx of cooperators more likely to trigger an institution to be formed. This is striking because the system is not biased in any particular direction and despite this still favours cooperation i.e., the institution being formed. This means that solely taking into account this fact of being a bounded agent in a noisy world increases the chances of an institution being formed.

The effect is more extreme for a larger range of noise, suggesting that a more diverse range of perceptions, due to the numerous factors affecting perception, further increases cooperation (Fig.4c).

Why does unbiased noise lead to a biased outcome? Despite the noise's unbiased nature, there is an asymmetry in the system. This is due to the proportional nature of the bias. To create a bias for each agent we either make them underestimate the expected cost of being punished by taking the product of a value less than 1 or above 1 to overestimate. Underestimating bias when the number of monitors is near zero does not change the estimate much, but overestimating has a larger net effect. This effect does not exist on the other side of the range (when there are many monitors) since there isn't a limit on how high one can overestimate.

5.4 Absolute noise (bounded uncertainty)

In the above section, we described the effect of a noisy proportional bias where the utility of an action is either over or under-valued in a proportional manner, but we could also implement this as absolute noisy perception, where an absolute value is added or subtracted to the utility of an action to produce a noisy expectation i.e. it can be in a range -8:8, -16:+16 which is added to the C_f estimate.

Under absolute noisy bias, we get a qualitatively different effect emerging (Fig.5), instead, it seems that the shifts in strategy space just become more

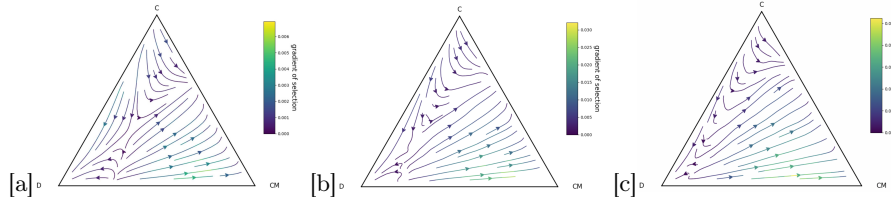


Fig. 4. (a) Control (b) noisy perception in range $[0.25:4]$ centred at 1 (c) noisy perception in range $[0.125:8]$ centred at 1. Note that due to the stochastic nature of the moran process used to model these dynamics, we cannot derive stable and unstable attractors (denoted by black and white dots) as we could with the deterministic replicator equation used in figures 1-3.

uncertain. We see that regions are more mixed and certain regions in strategy space that are far away from each other are more likely to be linked. This kind of noisy bias can help solve the bootstrapping problem since it can result in sudden shifts from non cooperative to cooperative equilibria. However, there would need to be a mechanism to reduce noise in estimates once a cooperative state is reached, to prevent the institution falling apart due to noise.

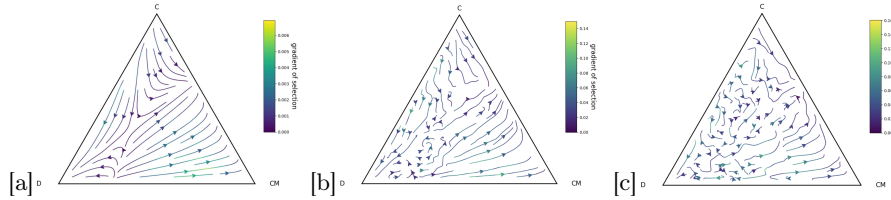


Fig. 5. (a) Control (b) absolute noisy perception range $[-8:8]$ (c) absolute noisy perception range $[-16:16]$. Note that due to the stochastic nature of the moran process used to model these dynamics, we cannot derive stable and unstable attractors (denoted by black and white dots) as we could with the deterministic replicator equation used in figures 1-3.

6 Discussion

6.1 Summary

Even when an institution is favourable to the individuals that constitute it, there is still the issue of meeting the critical mass of people joining it to incentivise its creation: no one wants to join and then be let down by no one else joining. We call this the institution bootstrapping problem.

To overcome this problem, we ask, what if the perceptual make-up of individuals makes them more likely to join an institution? And can merely believing

an institution exists, be enough to create it? Incorporating the psychological literature of subjective probability estimation and cognitive biases, we relax the assumption of perfect perception in game theory to investigate this. We show that a coarse bias in the subjective probability of being punished can make an institution either more likely or unlikely to form, this depends on whether agents have a bias that makes them either over or underestimate the subjective probability of being punished.

Further, we incorporate empirical work from psychophysics [38] which shows that individuals distort probabilities at the extremes of the range, biasing subjective probability of events when they are either highly unlikely or highly likely. We incorporate this and show that it also influences the threshold of monitors needed to form an institution in a similar fashion to a coarse grained bias. However, as described in the results, empirical work is unclear concerning how these biases work in social settings [38, 29].

We also find, that surprisingly, having a noisy perception across the population (no bias at group level), which is equivalent to individuals being biased by their bounded position in social and physical world, lowers the threshold of individuals needed to start an institution. Therefore this bounded noise is beneficial for forming a cooperative institution. Just the bare bone limitations of being a bounded agent embedded in a social world (hence noise) increase the chance of institutions being formed and make them more robust. This is if the noisy perception is proportional, however, if it is absolute, it generally makes the system more likely to shift states in an unpredictable manner (either to form an institution or not). This can be useful, however, there would need to be a way to eliminate this noisy absolute bias once the system is in a cooperative state i.e. agents would need to agree on a ground truth when conditions are favourable for cooperation.

Moreover, by relaxing game theory’s assumptions, which describes phenomena with group level State Variable Models (SVM) [34] e.g. average expected cost of free riding. We show, that even if you have a noisy uniform perception that cancels out and is therefore equivalent to the game theoretic SVM, the variance of perceptions can still change the outcome of the analysis.

6.2 Implications of this finding

Our findings suggest that: making agents have a more accurate estimate of the world, although it seems an intuitively good idea, may make them less likely to form an institution. Uncertainty and boundedness are therefore important to have. This rings true with other findings in other contexts which also stipulate that not having perfect information may be better for collective outcomes [31, 37, 6].

Further, we assume in this paper that the institution is good and benefits all agents who join it. But there may also be the case where agents may have been fooled into joining an institution that begets some but hurts others e.g. a corrupt institution where monitors take the collective resource for themselves. In this case, how would agents disobey an unjust institution [5]? In these cases, you

may want perfect/better information and for agents to agree with one another in order to rise up against the institution. This could be done if bounded agents come together by combining their estimates and taking the mean/median, as done in collective intelligence studies [16], and then, if they all use this group estimate in their decision-making, overcome a corrupt institution.

Furthermore, we can connect our model to collective action models [12] e.g. voting, taking political action, or protesting. If we interpret the cost of contributing as the cost that may be involved in collective action e.g. group effort, personal sacrifice e.t.c. then noise (bounded uncertainty) makes the institution (collective action) more likely, and having too much information (or too accurate an estimate) makes things hopeless, and agents stop trying to join the collective action. This corroborates well with Herbert Gintis' argument [12] that rational actors with perfect information are less likely to act since they perceive their actions will make no difference in large groups. He then asks: why do we see collective action even when it is irrational? Gintis then argues, drawing from the anthropological literature, that humans have an evolved bounded psychology that is tuned for small groups which changes their estimates of making a difference at the collective level and, therefore, explains the prevalence of collective action in human societies.

In conclusion, as shown in this paper, cognitive limitations of agents may actually be a benefit rather than a hindrance. Our general model can thus serve as an initial template to assess the impact of cognitive biases and uncertainty in more specific cases of institutions in socio-technical systems e.g. shared online servers, data sharing systems, and neighbourhood power cooperatives [32, 23, 22].

6.3 Our results contradict the conventional wisdom that noise is detrimental to cooperation

Early influential game theoretic and agent based models showed noise is deleterious for cooperation and leads to a breakdown of the tit-for-tat mechanism as it leads to rounds of mutual defection [3, 2].

Further, in an empirical study show that Salahshour et al [30] show that stochastic punishment (in this case the punishment factor - how much the fine is multiplied, which would be equivalent to having a noisy s in our model) reduces group contributions to the game and encourages more antisocial punishment (punishment of cooperators instead of defectors, which in turn disincentivises cooperation).

This discordance in our current paper about the positive collective benefits of noise may be explained by the fact that in our model, it is not noise in the punishment, but noise in the perception of the expected punishment that seems to engender cooperation, which is an important difference. These contextual factors make it hard to make general statements on the effects of noise on cooperation, one way or another. The nuance of uncertainty and noise and its non-trivial impact on norms and institutional emergence should be considered in the design of normative multiagent systems [32, 23, 22].

6.4 Limitations and extensions

So far we only have only one expected value in our model (expected punishment). However, there are other world variables that agents could have biased or noisy perceptions of e.g. cost of monitoring/cost of contribution. This can be altered in an agent based model, where it is possible to dissociate perceived and real payoff as in [21]. This allows the exploration of what factors in the population e.g. scarcity of resources lead to which biases/noise to evolve e.g. loss aversion bias [19, 1]

Further, although using EGT is a rigorous way to study the value of choices in a social environment and hence the direction in strategy space the system will move in. It stays agnostic to complex behaviours and interactions between individual and population-level dynamics. Agent-based models have the ability to capture complex behaviours and interactions in executable form, and to explore emergent phenomena simply by "running" variants of the model [8].

For example, a candidate mechanism for why uncertainty may encourage the joining of institutions: an initial proportion of individuals with noisy perception will be biased towards overestimating the probability of being caught, they will change their behaviour, which will in turn cause others who are less biased to do the same since there will now be more members in the institution, this process repeats instantiating a feedback loop that allows agents to bootstrap to cooperation.

Acknowledgments. Stavros Anagnou was supported by a PhD studentship from the University of Hertfordshire and a MITACS Globalink Research Award grant. This research was made possible, in part, thanks to the Canada Research Chairs program. We would also like to thank the anonymous reviewers for their feedback on the manuscript and Elias Fernández Domingos for advice regarding the EGTtools library.

Disclosure of Interests. The authors have no competing interests to declare that are relevant to the content of this article.

References

1. Anagnou, S., Polani, D., Salge, C.: The effect of noise on the emergence of continuous norms and its evolutionary dynamics. In: Proceedings of the ALIFE 2023: Ghost in the Machine: Proceedings of the 2023 Artificial Life Conference. p. 125. MIT Press. https://doi.org/https://doi.org/10.1162/isal_a_00588
2. Axelrod, R., Dion, D.: The further evolution of cooperation **242**(4884), 1385–1390. <https://doi.org/10.1126/science.242.4884.1385>, <https://www.science.org/doi/10.1126/science.242.4884.1385>
3. Bendor, J., Kramer, R.M., Stout, S.: When in doubt...: Cooperation in a noisy prisoner's dilemma **35**(4), 691–719. <https://doi.org/10.1177/0022002791035004007>, <https://journals.sagepub.com/doi/10.1177/0022002791035004007>
4. Boyd, R., Richerson, P.J.: Culture and the evolutionary process. University of Chicago Press, paperback edn.

5. Burth Kurka, D., Pitt, J., Lewis, P.R., Patelli, A., Ekart, A.: Disobedience as a mechanism of change. In: 2018 IEEE 12th International Conference on Self-Adaptive and Self-Organizing Systems (SASO). pp. 1–10. IEEE. <https://doi.org/10.1109/SASO.2018.00011>, <https://ieeexplore.ieee.org/document/8613714/>
6. Davies, A.P., Watson, R.A., Mills, R., Buckley, C.L., Noble, J.: “if you can’t be with the one you love, love the one you’re with”: How individual habituation of agent interactions improves global utility **17**(3), 167–181. https://doi.org/10.1162/artl_a_00030, <https://direct.mit.edu/artl/article/17/3/167-181/2687>
7. Duguid, S., Wyman, E., Bullinger, A.F., Herfurth-Majstorovic, K., Tomasello, M.: Coordination strategies of chimpanzees and human children in a stag hunt game **281**(1796), 20141973. <https://doi.org/10.1098/rspb.2014.1973>, <https://royalsocietypublishing.org/doi/10.1098/rspb.2014.1973>
8. Epstein, J.M., Axtell, R.: Growing artificial societies: social science from the bottom up. Complex adaptive systems, Brookings Institution Press
9. Fernández Domingos, E., Santos, F.C., Lenaerts, T.: EGTtools: Evolutionary game dynamics in python **26**(4), 106419. <https://doi.org/10.1016/j.isci.2023.106419>, <https://linkinghub.elsevier.com/retrieve/pii/S2589004223004960>
10. Frey, S., Sumner, R.W.: Emergence of integrated institutions in a large population of self-governing communities **14**(7), e0216335. <https://doi.org/10.1371/journal.pone.0216335>, <https://dx.plos.org/10.1371/journal.pone.0216335>
11. Galton, F.: Vox populi **75**(1949), 450–451. <https://doi.org/10.1038/075450a0>, <https://www.nature.com/articles/075450a0>
12. Gintis, H.: Individuality and Entanglement: The Moral and Material Bases of Social Life. Princeton University Press. <https://doi.org/10.2307/j.ctvc779cx>, <http://www.jstor.org/stable/10.2307/j.ctvc779cx>
13. Hertwig, R., Herzog, S.M.: Fast and frugal heuristics: Tools of social rationality **27**(5), 661–698. <https://doi.org/10.1521/soco.2009.27.5.661>, <http://guilfordjournals.com/doi/10.1521/soco.2009.27.5.661>
14. Johnson, D.D.P., Fowler, J.H.: The evolution of overconfidence **477**(7364), 317–320. <https://doi.org/10.1038/nature10384>, <https://www.nature.com/articles/nature10384>
15. Kahneman, D., Sibony, O., Sunstein, C.R.: Noise: a flaw in human judgment. William Collins
16. Kao, A.B., Berdahl, A.M., Hartnett, A.T., Lutz, M.J., Bak-Coleman, J.B., Ioannou, C.C., Giam, X., Couzin, I.D.: Counter-acting estimation bias and social influence to improve the wisdom of crowds **15**(141), 20180130. <https://doi.org/10.1098/rsif.2018.0130>, <https://royalsocietypublishing.org/doi/10.1098/rsif.2018.0130>
17. Krause, S., James, R., Faria, J.J., Ruxton, G.D., Krause, J.: Swarm intelligence in humans: diversity can trump ability **81**(5), 941–948. <https://doi.org/10.1016/j.anbehav.2010.12.018>, <https://linkinghub.elsevier.com/retrieve/pii/S0003347210005221>
18. Lewis, P.R., Marsh, S.: What is it like to trust a rock? a functionalist perspective on trust and trustworthiness in artificial intelligence **72**, 33–49. <https://doi.org/10.1016/j.cogsys.2021.11.001>, <https://linkinghub.elsevier.com/retrieve/pii/S1389041721000814>

19. McNamara, J., Houston, A.: Risk-sensitive foraging: A review of the theory **54**(2), 355–378. [https://doi.org/10.1016/S0092-8240\(05\)80031-X](https://doi.org/10.1016/S0092-8240(05)80031-X), [http://link.springer.com/10.1016/S0092-8240\(05\)80031-X](http://link.springer.com/10.1016/S0092-8240(05)80031-X)
20. McNamara, J.M.: Game theory in biology: Moving beyond functional accounts **199**(2), 179–193. <https://doi.org/10.1086/717429>, <https://www.journals.uchicago.edu/doi/10.1086/717429>
21. McNamara, J.M., Houston, A.I., Leimar, O.: Learning, exploitation and bias in games **16**(2), e0246588. <https://doi.org/10.1371/journal.pone.0246588>, <https://dx.plos.org/10.1371/journal.pone.0246588>
22. Morris-Martin, A., De Vos, M., Padget, J.: Norm emergence in multiagent systems: a viewpoint paper **33**(6), 706–749. <https://doi.org/10.1007/s10458-019-09422-0>, <http://link.springer.com/10.1007/s10458-019-09422-0>
23. Nardin, L.G., Balke-Visser, T., Ajmeri, N., Kalia, A.K., Sichman, J.S., Singh, M.P.: Classifying sanctions and designing a conceptual sanctioning process model for socio-technical systems **31**(2), 142–166. <https://doi.org/10.1017/S0269888916000023>, <https://www.cambridge.org/core/product/identifier/S0269888916000023/type/journal-article>
24. North, D.C.: Institutions, Institutional Change and Economic Performance. Cambridge University Press, 1 edn. <https://doi.org/10.1017/CBO9780511808678>, <https://www.cambridge.org/core/product/identifier/9780511808678/type/book>
25. Ostrom, E.: Governing the commons: the evolution of institutions for collective action. The Political economy of institutions and decisions, Cambridge University Press
26. Pachur, T., Hertwig, R., Steinmann, F.: How do people judge risks: Availability heuristic, affect heuristic, or both? **18**(3), 314–330. <https://doi.org/10.1037/a0028279>, <https://doi.apa.org/doi/10.1037/a0028279>
27. Pitt, J., Schaumeier, J., Artikis, A.: Axiomatization of socio-economic principles for self-organizing institutions: Concepts, experiments and challenges **7**(4), 1–39. <https://doi.org/10.1145/2382570.2382575>, <https://dl.acm.org/doi/10.1145/2382570.2382575>
28. Powers, S.T., Ekárt, A., Lewis, P.R.: Modelling enduring institutions: The complementarity of evolutionary and agent-based approaches **52**, 67–81. <https://doi.org/10.1016/j.cogsys.2018.04.012>, <https://linkinghub.elsevier.com/retrieve/pii/S1389041718301268>
29. Rilling, J.K., Sanfey, A.G., Aronson, J.A., Nystrom, L.E., Cohen, J.D.: The neural correlates of theory of mind within interpersonal interactions **22**(4), 1694–1703. <https://doi.org/10.1016/j.neuroimage.2004.04.015>, <https://linkinghub.elsevier.com/retrieve/pii/S1053811904002241>
30. Salahshour, M., Oberhauser, V., Smerlak, M.: The cost of noise: Stochastic punishment falls short of sustaining cooperation in social dilemma experiments **17**(3), e0263028. <https://doi.org/10.1371/journal.pone.0263028>, <https://dx.plos.org/10.1371/journal.pone.0263028>
31. Salge, C., Polani, D.: Don’t believe everything you hear: Preserving relevant information by discarding social information. In: Artificial Life 14: Proceedings of the Fourteenth International Conference on the Synthesis and Simulation of Living Systems. pp. 837–844. The MIT Press. <https://doi.org/10.7551/978-0-262-32621-6-ch137>, <https://www.mitpressjournals.org/doi/abs/10.1162/978-0-262-32621-6-ch137>

32. Savarimuthu, B.T.R., Cranefield, S.: Norm creation, spreading and emergence: A survey of simulation models of norms in multi-agent systems **7**(1), 21–54. <https://doi.org/10.3233/MGS-2011-0167>, <https://journals.sagepub.com/doi/full/10.3233/MGS-2011-0167>
33. Simon, H.A.: A behavioral model of rational choice **69**(1), 99. <https://doi.org/10.2307/1884852>, <https://academic.oup.com/qje/article-lookup/doi/10.2307/1884852>
34. Smith, J.M.: Evolution and the Theory of Games. Cambridge University Press, 1 edn. <https://doi.org/10.1017/CBO9780511806292>, <https://www.cambridge.org/core/product/identifier/9780511806292/type/book>
35. Sundh, J.: Human behavior in the context of low-probability high-impact events **11**(1), 902. <https://doi.org/10.1057/s41599-024-03403-9>, <https://www.nature.com/articles/s41599-024-03403-9>
36. Tversky, A., Kahneman, D.: Advances in prospect theory: Cumulative representation of uncertainty **5**(4), 297–323. <https://doi.org/10.1007/BF00122574>, <http://link.springer.com/10.1007/BF00122574>
37. Vogrin, M., Wood, G., Schmickl, T.: Confirmation bias as a mechanism to focus attention enhances signal detection **26**(1), 2. <https://doi.org/10.18564/jasss.4954>, <https://www.jasss.org/26/1/2.html>
38. Zhang, H., Ren, X., Maloney, L.T.: The bounded rationality of probability distortion **117**(36), 22024–22034. <https://doi.org/10.1073/pnas.1922401117>, <https://pnas.org/doi/full/10.1073/pnas.1922401117>