

# Evolutionary Machine Ethics

The Anh Han<sup>1</sup> & Luís Moniz Pereira<sup>2</sup>

School of Computing, Media and the Arts,  
Teesside University, TS1 3BX, UK.  
Email: [T.Han@tees.ac.uk](mailto:T.Han@tees.ac.uk)

NOVA Laboratory for Computer Science and Informatics (NOVA LINCS)  
Departamento de Informática, Faculdade de Ciências e Tecnologia  
Universidade Nova de Lisboa, 2829-516 Caparica, Portugal.  
Email: [imp@fct.unl.pt](mailto:imp@fct.unl.pt)

**Abstract.** Machine ethics is a sprouting interdisciplinary field of enquiry arising from the need of imbuing autonomous agents with some capacity for moral decision-making. Its overall results are not only important for equipping agents with a capacity for moral judgment, but also for helping better understand morality, through the creation and testing of computational models of ethics theories. Computer models have become well defined, eminently observable in their dynamics, and can be transformed incrementally in expeditious ways. We address, in work reported and surveyed here, the emergence and evolution of cooperation in the collective realm. We discuss how our own research with Evolutionary Game Theory (EGT) modelling and experimentation leads to important insights for machine ethics, such as the design of moral machines, multi-agent systems, and contractual algorithms, plus their potential application in human settings too.

**Keywords:** Machine Ethics; Evolutionary Game Theory; Intention Recognition; Commitment; Guilt.

## 1. Introduction

Some of our previous research (Pereira & Saptawijaya, 2011; Han, Saptawijaya, & Pereira, 2012; Pereira & Saptawijaya, 2015, 2016a, 2016b, 2017; Saptawijaya & Pereira, 2015, 2016a, 2018) has focused on using logic programming techniques to computational modelling of morality sans emotions. In the realm of the individual, we have addressed questions of permissibility and the dual-process of moral judgments by framing together ingredients that are essential to moral agency: abduction, integrity constraints, preferences, argumentation, counterfactuals, and updates. Computation over these ingredients has become our vehicle for modelling the dynamics of moral cognition within a single agent, without addressing the cultural dimension (Prinz, 2016), because it is still absent in machines.

In the collective realm, we have reported on computational moral emergence (Han et al., 2015a), again sans emotions, using techniques from Evolutionary Game Theory (EGT). We have shown that the introduction of cognitive abilities, like intention recognition, commitment, revenge, apology, forgiveness and guilt, reinforce the emergence of cooperation in diverse populations, comparatively to their absence, by way of EGT models. This evolutionary collective realm will be the one surveyed here, with the pointers to our specialized publications to be indicated below.

In studies of human morality, these distinct but interconnected realms – one stressing above all individual cognition, deliberation, and behaviour; the other stressing collective morals and how they have emerged with evolution – seem separate but are synchronously evinced (Pereira & Saptawijaya, 2015, 2016b). There are issues concerned with how to bridge the two realms also addressed in (Gaspar, 2016). Our account affords plenty of room for an evolutionary

phylogenetic emergence of morality, as illustrated in the sequel, thereby supplementing the limitations of focusing just on ontogeny. The bridging issues concern individual cognitive abilities and their deployment in the population. Namely the ability of recognising intention in another, taking even into account how others recognize our intention; the abilities of requesting commitment, and of accepting or declining to commit; the abilities to adaptively apply complementary mechanisms; those of monitoring group participation and delegate this process to an external party; those of cooperating or defecting; plus those of apologising, be it fostered by guilt, and of taking revenge or forgiving.

This chapter relies mainly on our collective realm research, and considers the modelling of distinct co-present strategies of cooperative and uncooperative behaviour. Such driving strategies are associated with moral “emotions” that motivate moral discernment and substantiate ethical norms, leading to improved general conviviality on occasion, or not. To wit, we can model moral agency without explicitly representing embodied emotions as we know them. Rather, such software-instantiated “emotions” are modelled as (un)conscious heuristics empowered in complex evolutionary games.

In the next two sections, starting with the ground breaking work of Alan Turing, functionalism is employed to scaffold a philosophical perspective on emotions and morality. The further seven sections after those review materials from our EGT-based research in support of this perspective. This work has substantiated the philosophical viewpoint through an admixture of intention recognition, commitment, revenge, apology, forgiveness, and guilt. The final section conjectures further on guilt, and its relationship with counterfactual reasoning, as a next natural step in our research program.

## **2. Turing is Among Us**

Turing's relevance arises from the timelessness of the issues he tackled, and the innovative light he shed upon them (Pereira, 2012a). He first defined the algorithmic limits of computability, via an *effective* well-specified mechanism, and showed the generality of his definition by proving its equivalence to other general, but less algorithmic and non-mechanical, more abstract formulations of computability. His originality lies on the essential simplicity of the mechanism invoked – the now dubbed Turing Machines (or programs), which he called A-Machines – and the proof of existence of a Universal A-Machine (i.e. the digital computer, known in academia as the Universal Turing Machine), which can simulate any other A-Machine, that is, execute any program.

Interestingly, he raised the issue of whether human beings are a measure for his “machines”, and, in mechanizing human cognition Turing implicitly introduced the modern perspective since known as “functionalism”. According to this paradigm, what counts is the realisation of function, independently of the hardware embodying it. Such “multiple realisation” is afforded by the very simplicity of his devised mechanism, relying solely on the manipulation of discrete information, where data and instructions are both represented just with symbols. The twain are stored in memory, instructions doubling as data and as rules for acting – the stored program. To this day, no one has invented a computational mechanical process with such general properties, which cannot be theoretically approximated with arbitrary precision by some A-Machine, where any interactions with the world outside are captured by Turing's innovative concept and definition of “oracle” – the very word employed by him for the purpose –, as a means to interrogate that world by posing queries to one or more outside oracles. This concept of oracle is regularly taught in computer science today, namely in the essential study of computation complexity, though not every student knows it came from Turing. In the midst of a computation a query may be posed to an outside oracle about the satisfaction of some truth, and the computation continued once an answer obtained, rather than the computer testing for an answer in a possibly infinite set of them.

Turing further claimed that his machines could simulate the effect of *any* activity of the mind, not just a mind engaged upon a “definite method of proceeding” or algorithm. He was clear that discrete state machines included those with learning or self-organising abilities, and stressed that these still fall within the scope of the computable. Turing drew attention to the apparent conflict between self-organisation and the definition of A-Machines as having fixed tables of behaviour, but sketched a proof that self-modifying machines are still definable by an unchanged instruction set (Hodges 1997; McDermott 2001).

The promise of this approach in studies of morality is that it represents a universal functionalism, the terms of which enable the bringing together of the ghosts in the several embodied machines (silicon-based, biological, extra-terrestrial or otherwise), to promote their symbiotic epistemic co-evolution, as they undertake moral action within a common moral theatre.

### **3. Functionalism and Emergence**

The principle of the distinction between software and hardware appears clear-cut with the advent of the digital computer and its conceptual precursor, the Universal Turing Machine. The diversity of technologies employed to achieve the same function, confirms it ever since the first computers. One program is executable in physically different machines, precisely because the details of its execution below an ascertainable level of analysis are irrelevant, as long as an identical result at the level of discourse is produced. That said, however, the distinction between hardware and software is not so clear as it might seem. Hardware is not necessarily represented by things physical but rather by what, at some level of analysis, is considered fixed, given, and whose analysis or non-analysability is irrelevant for the purpose at hand. Historically, in the first computers, that level coincided with that of the physical parts of the machine. Subsequently, especially due to rapidly increasing computing power, “hardware” has become increasingly “soft”, with the physical basis for the hardware/software distinction finally blurred by the concept of the “abstract machine”: a fixed collection of mathematically defined instructions supporting a set of software functions, independently of the particular physical processes underlying the implementation of the abstract machine, that is, realising it.

Hence, “multiple realisation” stands for the thesis that a mental state can be “realised” or “implemented” by different physical states. Beings with different physical constitutions can thus be in the same mental state, and from these common grounds can cooperate, acting in mutual support (or not). According to classical functionalism, multiple realisation implies that psychology is autonomous: in other words, biological facts about the brain are irrelevant (Boden, 2008). Whether physical descriptions of the events subsumed by psychological generalisations have anything in common is irrelevant to the truth of the generalisations, to their interestingness, to their degree of confirmation, or, indeed, to any of their epistemological important properties (Fodor 1974).

Functionalism has continued to flourish, being developed into numerous versions by thinkers as diverse as David Marr, Daniel Dennett, Jerry Fodor, and David Lewis (Fodor 1974; Dennett 2005). It helped lay the foundations for modern cognitive science, being the dominant theory of mind in philosophy today. In the latter part of the 20th and early 21st centuries, functionalism stood as the dominant theory of mental states. It takes mental states out of the realm of the “private” or subjective, and gives them status as entities open to scientific investigation. Functionalism's characterisation of mental states in terms of their roles in the production of behaviour grants them the causal efficacy that common sense takes them to have. In permitting mental states to be multiply realised, functionalism offers an account of mental states compatible with materialism, without limiting the class of minds to creatures with brains like ours (Levin 2010).

Biological evolution is characterized by a set of highly braided processes, which produce a kind of extraordinarily complex combinatorial innovation. A generic term frequently used to describe

this vast category of spontaneous, and weakly predictable, order-generating processes, is “emergence”. This term became a sort of signal to refer to the paradigms of research sensitive to systemic factors. Complex dynamic systems can spontaneously assume patterns of ordered behaviours not previously imaginable from the properties of their constitutive elements or from their interaction patterns. There is unpredictability in self-organising phenomena – preferably called “evolutionary” (Turing 1950) – with considerably variable levels of complexity, where “complexity” refers to the emergence of collective properties in systems with many interdependent components. These components can be atoms or macromolecules in physical or biological contexts, and people, machines or organisations in socioeconomic contexts.

What does emerge? The answer is not something defined physically but rather something like a shape, pattern, or function. The concept of emergence is applicable to phenomena in which the relational properties predominate over the properties of the compositional elements in the determination of the ensemble’s characteristics. Emergence processes are due to starting configurations and interaction topologies, not intrinsic to the components themselves (Deacon 2003). This functionalism is, almost by definition, anti substance-essence, anti vital-principle, anti monopoly of *qualia*.

Building intelligent machines may seek a partial understanding of the emergence of higher-level properties, like morality. Here, functionalism affirms the salience of the results of this work in assessing, for example, human morality. Again, functionalism holds that the material substrate is not of the essence, and that it suffices to realise equivalent functionality albeit by way of a different material vehicle. Moreover, distinct roads to the same behaviour may be had, thereby adding to our understanding of what, say, “general intelligence” or “mind” means. Thus, on our estimation, the most fruitful inquires into the nature of “mind” or “general intelligence” will certainly include the use of Artificial Intelligence aided in time by the embryonic field of artificial emotions, qua strategies, to simulate complex mental operations, as already foreseen (Turing 1950).

#### **4. Evolutionary Game Theory**

Game theory was first developed in the 1940’s, and the first work on the subject was *Theory of Games and Economic Behavior* by the mathematician John von Neumann (1903-1957) and the economist Oskar Morgenstern (1902-1977), (Neumann & Morgenstern, 1944). At the time it was directed at the economy, but it was subsequently applied to the Cold War, as the outcome of issues raised by the use of the atomic bomb and the subtle means of bluffing. When some such situation gets complicated, there is need to resort to sophisticated mathematical tools — and computer simulations — to deal with equations that cannot otherwise be solved.

The games theme is as complex as it is interesting and filled with diverse niches. We can envisage genes and memes (“cultural genes”), and their mutual combinations, as ongoing strategies in the game of evolution, raising issues and posing questions related to survival and winning. We can envisage too the combinatorial evolution of such strategies, and their possible mutations according to diverse conditions, which conditions can either be other game partners or the game board rules (or Nature's) own circumstances. The notion of game includes uncertainty, and whenever there is uncertainty there has to be some attending strategy, spelling the moves one makes with given probability. When there is co-presence of evolving strategies from several partners, along with the idea of game payoff, we are dealing with the notion of evolutionary game, which can be examined and studied in an abstract and mathematical manner.

The same way we have genetic strategies for reproduction, all of our lives are filled with cultural, or civilizing, strategies. And, in a general way, we can view and scry our species through such lenses, still without undervaluing other remaining perspectives, equally important.

There are zero sum games and non-zero sum games. The zero sum ones are those that, by their rules, some players win, some players loose. In Nature’s evolution, conditions are those

of non-zero sum — all can win or all can lose. Robert Wright (2001) analyses the evolution of culture and civilization with the underlying idea that, in Nature, non-zero sum games are possible, wherefore a general gain may be obtained through cooperation, thereby leading to illuminated altruism.

Sometimes, co-present strategies tend to achieve a tactical equilibrium. Take the hunter/prey relationship: neither the hunter wants to fully exterminate the prey, nor the latter can multiply indefinitely because that would exhaust the environment's resources. Some of these studies are used by Economics to understand what might be the overall result from the sum of interactions amongst the several game partners.

It is relevant to take into consideration if the game takes place only once with a given partner, or whether the same partner may be encountered on other occasions; how much recall does one have of playing with that partner; and whether the possibility of refusing a partner is allowed. Let us take a more detailed look at each of these situations in turn. We begin with the famous Prisoners Dilemma (PD), typical of the paradox of altruism. There are two prisoners, A and B, with charges on them. Either of them can denounce the other, or confess, or remain silent.

	Prisoner <b>B</b> – silence	Prisoner <b>B</b> – confession
Prisoner <b>A</b> – silence	<b>6 years in jail for each</b>	<b>A = 10 years in jail</b> <b>B = 2 years in jail</b>
Prisoner <b>A</b> – confession	<b>A = 2 years in jail</b> <b>B = 10 years in jail</b>	<b>8 years in jail for each</b>

Consider the above 2x2 payoff matrix where the lines correspond to the behaviour of A (to remain silent or confess), and the columns correspond to the behaviour of B (to remain silent or confess). At the intersection of B's «confess» column with A's «confess» row, both receive a jail sentence of 8 years. If A confesses and B does not, A will only get a 2-year sentence, whereas B gets 10 years, and vice-versa. There is an incentive for any of them to confess in order to reduce their own jail sentence. This way, it would eventually be advantageous for them not to remain silent. If one of them defects by confessing, but not the other, he will only stay in jail for 2 years whereas the other will be there for 10 years. But if both confess they will be sentenced to 8 years each. The temptation to confess is great, but so is the inherent risk, because, after all, they would mutually benefit from remaining silent, getting a 6-year sentence each in that case.

The prisoners know the rules of the game; they just do not know how the other player will act. It is advantageous for them to remain silent, but they do not know if the other one will confess. As long as one of them confesses, the silent other will be sentenced to 10 years in jail. A dilemma thus arises: it is good to remain silent, but there is the risk the other one will defect; and the one who does it faster will take the greater gain. In the worst-case scenario, both get an 8-year sentence — nobody will take the risk. This is a classic game, one where both players have the tendency to confess — and not benefit from what could be a mutual advantage, but one that they cannot assuredly profit from. Firstly, they do not have the opportunity to talk to one another; secondly, because even if they did, they would still risk being betrayed by the other. They have no joint solution in the sense that A and B could ever choose what is best for both, where there would be an assured increased advantage for the two.

All turns more complicated when one imagines A and B playing this game many times in succession, taking into account their experience of previous mutual behaviour in their past. In this case they can go on building mutual trust or distrust. If one betrayed the other once, the

betrayed one's reaction will be vengeance, or simply intolerance, in some future opportunity. Let us now visualize a situation with multiple players and ask ourselves which will be, along time, the best of all possible strategies — by running a computer simulation. Of course one thing is to presuppose any one strategy can always match with any other, which is the base assumption, and then to move on to a situation where one wants to match only with certain players. Through these more realistic situations one begins to develop a game theory where social structure is included inside it.

Instead of letting a strategy evolve by choosing to copy those who win the most, one can alternatively let those who win the most to be those that reproduce the most, that is, they make more copies of themselves proportionately to the others, all the while keeping a bounded size for the whole population (since overall resources are not unlimited), through a random elimination of individuals. This other option can be adopted because those who lose more (or win less) are eliminated by virtue of their reduced number of copies, and also because only those who win more than some threshold are allowed to reproduce (reproduction is costly). The intent of this interpretation is that, throughout the game, strategies want to take over resources and occupy vital space in the population. Winning means having more energy to reproduce, while losing means not being able to persist with one's genetic/memetic continuity.

The evolutionary question that arises then is whether everyone can at length benefit more if they cooperate more. Which question hinges on how to prevent free-riders who want to gain more without having to incur in the expenses of cooperation. The evolution of any collective species clashes against this problem of balancing cooperation with opportunism. It is a strong theme in Evolutionary Psychology (Pereira 2012b), and one for which we can devise mathematical models and employ computers to perform both analytical computations, as well as long and repetitive simulations of the joint evolution of behavioural strategies in co-presence, typically done via mathematical games' implementation mixing competitive and cooperative situations, and providing mutation in strategies in order to detect focus points of long-term evolution stability.

## **5. Learning to recognise intentions and committing can resolve cooperation dilemmas**

Few problems have motivated the amalgamation of so many seemingly unrelated research fields as has the evolution of cooperation (Nowak, 2006). Several mechanisms have been identified as catalysers of cooperative behaviour; see for example surveys by Nowak (2006) and Sigmund (2010). Yet these studies, mostly grounded on evolutionary dynamics and game theory, have neglected the important role played by intention recognition (Han and Pereira, 2013c) in behavioural evolution. In our work (Han et al., 2011, 2012a, 2012b; Han, 2013), we explicitly studied the role of intention recognition in the evolution of cooperative behaviour. The results indicate that intention recognisers prevail against the most successful strategies in the context of the Iterated Prisoner's Dilemma (IPD) (e.g. win-stay-lose-shift, and tit-for-tat like strategies), and promote a significantly higher level of cooperation, even in the presence of noise, plus the reduction of fitness associated with the cognitive costs of performing intention recognition. Our approach offers new insights into the complexity of – as well as enhanced appreciation for the elegance of – behavioural evolution when driven by elementary forms of cognition and learning ability.

It is important to note that intention recognition techniques have been studied actively in AI for several decades (Charniak and Goldman, 1993; Sadri, 2011; Han and Pereira, 2013a, 2013b, 2013c), with various applications such as for improving human-computer interactions, assistive living, moral reasoning, and team work (Pereira and Han, 2011a, 2011b; Roy et al., 2007; Han et al., 2012d; Heinze, 2003; Han and Pereira, 2013b). Intentionality has been also shown to play a crucial role in making moral judgments, e.g. as captured in the Doctrines of Double and of Triple Effect (Hauser, 2006; Mikhail, 2007). Therefore, our results, both analytically and through extensive agent-based simulations, provide important insights into designing of moral agents and machines that are capable recognising others' intentions and taking them into

account in their moral decision judgement. A clear implication is that, by virtue of such designs, moral agents in a society will be able to maintain high levels of cooperative behaviours.

Now, conventional wisdom suggests that clear agreements need to be made prior to any collaborative effort in order to avoid potential frustrations for the participants. We have shown (Han et al., 2013a) that this behaviour may actually have been shaped by natural selection, as argued in (Nesse, 2011). Our research demonstrates that reaching prior explicit agreement about the consequences of not honouring a deal provides a more effective road to facilitating cooperation than simply punishing bad behaviour after the fact, even when there is a cost associated to setting up the explicit agreement. Typically, when starting a new project in collaboration with someone else, it pays to establish up-front how strongly your partner is prepared to commit to it. To ascertain the commitment level one can ask for a pledge and stipulate precisely what will happen should the deal not be honoured.

In our study, EGT is used to show that when the cost of arranging commitments (for example, that of hiring a lawyer to make a contract) is justified with respect to the benefit of the joint endeavour (for instance buying a house), and that, when the compensation is set sufficiently high, commitment proposers become prevalent, thence leading to a significant level of cooperation. Commitment proposers can get rid of fake co-operators that agree to cooperate with them yet act differently, thus also avoiding interaction with the bad guys that only aim to exploit the efforts of the cooperative ones. Interestingly, we have shown that whenever the compensation cost reaches a certain threshold (roughly equal the sum of the cost of arrangement commitment plus the benefit of cooperation), no further improvement is achieved by increasing the compensation. This outcome implies that, for regulating legal contracts, it is not required to set extreme penalties for small issues, which might otherwise lead to undesirable side-effects, such as the unwillingness to commit due to the contracts figuring extreme penalties.

But what happens if the cost of arranging the commitments is too high compared to the benefit of cooperation? Would you make a legal contract for sharing a cake? Our results show that in that case those that free ride on the investment of others will “immorally” and inevitably benefit. Establishing costly agreements only makes sense for specific kinds of projects. Our study shows that insisting your partner jointly share in the cost of setting up a deal leads to even higher levels of cooperation, supporting the evolution of cooperation for a larger range of arrangement costs and compensations. This makes sense, as equal investment will ensure the credibility of the pledge by both partners. Agreements based on shared costs result in better friends.

More interestingly, our research (Han, et al., 2015a, Han, et al., 2015b) into the synergy of the two presented mechanisms, those of intention recognition and prior commitment, sheds new light on promoting cooperative behaviour. This work employs EGT methods in agent-based computer simulations to investigate mechanisms that underpin cooperation in differently composed societies. High levels of cooperation can be achieved if reliable agreements can be arranged. Formal commitments, such as contracts, promote cooperative social behaviour if they can be sufficiently enforced, and the costs and time to arrange them provide mutual benefit.

On the other hand, an ability to assess intention in others has been demonstrated to play a role in promoting the emergence of cooperation. Indeed, this ability to assess the intentions of others based on experience and observations facilitates cooperative behaviour without resort to formal commitments like contracts. To wit, our research found that the synergy between intention recognition and commitment strongly depends on the confidence and accuracy of the intention recognition capacity. To reach high levels of cooperation, commitments might be unavoidable whenever intentions cannot be assessed with sufficient confidence and accuracy. Otherwise, it is advantageous to wield solely intention recognition so as to avoid the costly arranging of commitments.

In short, it seems to us that intention recognition, and its use in the scope of commitment, is a foundational cornerstone where we should begin at, naturally followed by the capacity to establish and honour commitments, as a tool towards the successive construction of collective intentions and social organization (Searle, 1995, 2010). Moreover, given that arranging prior commitments is a way to reveal others' intentions (Cohen and Levesque, 1990), our sustained hope has been that the combination of these two complementary mechanisms provides useful implications for the design of moral machines that are capable of better intention prediction (e.g. when no prior information is available regarding the recognised agents) and intention-based moral judgements.

## **6. Combining commitment and costly punishment to prevent antisocial behaviour**

We have compared prior commitment with costly posterior punishment, a strategy that makes no prior agreements at all and simply punishes wrongdoers afterwards. Previous studies show that, by punishing bad behaviour strongly enough, cooperation can be promoted in a population of self-interested individuals, see e.g. (Fehr & Gächter, 2002; Han, 2016). Yet these studies also show that the punishment must sometimes be quite excessive in order to obtain significant levels of cooperation. Our own study shows that arranging prior agreements can significantly reduce the impact-to-cost ratio of punishment. Higher levels of cooperation can actually be attained by dint of lower levels of punishment.

More interestingly, through the observation that prior commitment and posterior punishment complement each other, nicely dealing with different types of defective behaviours, we investigated different ways in which these two strategies can be combined. First of all, in (Han and Lenaerts 2016), we have shown that a simple probabilistic combination of the two mechanisms can promote a higher level of cooperation rather than either commitment or punishment alone. It is based on the assessment that arranging prior commitment reduces the effect-to-cost ratio required by costly punishment to perform efficiently, particularly when the cost of arrangement is sufficiently low. While costly punishment can enable one to deal with commitment free-riders, i.e. those who can escape sanctioning when interacting with the commitment strategy simply by avoiding commitment. Our analytical and simulation results show that a combined strategy leads to substantial enhancement in terms of the level of cooperation. Notably, this level is most significant when the punishment cost is sufficiently large and the impact of punishment reaches a threshold. As such, our results have shown that the combined strategy can simultaneously overcome the weaknesses of both strategies.

We have studied another combination approach to exploiting the complementarities of the two mechanisms, in which they are now co-present in the population (Han 2016). Interestingly, it provides a novel solution to prevent antisocial punishment, the one where defectors can punish cooperators, a major challenge in the studies of the evolution of cooperation (Raihani and Bshary, 2015; Power et al., 2011). Namely we have shown, in the context of the one-shot PD, that, if in addition to using punishment the agents in a population can also propose cooperation agreements to their co-players prior to an interaction, then social punishment and cooperation can evolve together, even in the presence of said antisocial punishment. Antisocial punishers can be significantly restrained by commitment proposing agents since only those who dishonour a commitment deal can be enforced to pay compensation. On the other hand, since arranging a commitment deal is costly, its regime can be replaced by social punishers who do not have to pay this cost, while still maintaining cooperation among them. Our results have shown that when both strategic options of commitment and punishment are present, social punishment dominates a population with antisocial punishment players, leading to a significantly higher level of cooperation compared to the cases when either of the strategic options is absent. This is a notable observation since arranging prior commitments, by itself, is already a strong mechanism that can enforce a substantial level of cooperation. But by sacrificing via the extra cost of commitment for a punishment strategy that is vulnerable to antisocial behaviours and

defection, it then results in a significant improvement in terms of cooperation. That is, the commitment mechanism catalyses the emergence of social punishment and cooperation.

In short, our results provide novel insights for the design of autonomous and multi-agent systems comprised of moral agents that require cooperation amongst them in a competitive environment, especially when they are based on commitments or on sanctions to regulate agents' behaviours. We have shown for the first time that combining the two mechanisms can lead to a better strategy for cooperation promotion and, furthermore, prevent antisocial behaviours whilst simultaneously maximising the benefit of deploying an appropriate sanctioning system.

## **7. Commitments can resolve group cooperation dilemmas. On Avoidance, Restriction, Participation Monitoring, and Delegation**

Public goods, like food sharing and social health systems, may prosper when prior agreements to contribute are feasible and all participants commit to do so. Yet, free-riders may exploit such agreements (Han et al., 2013a), thus requiring committers to decide not to enact the public good whenever sufficient others are not attracted to committing. This decision removes all benefits from free-riders (non-contributors), but also from those who are wishing to establish the beneficial resource. In (Han et al., 2014) we show, in the framework of the one-shot Public Goods Game (PGG) and EGT, that implementing measures to delimit benefits to “immoral” free-riders, often leads to more favourable societal outcomes, especially in larger groups and in highly beneficial public goods situations, even if doing so incurs in new costs. PGG is the standard framework for studying emergence of cooperation within group interaction settings (Sigmund, 2010). In a PGG, players meet in groups of a fixed size, and all players can choose whether to cooperate and contribute to the public good or to defect without contributing to it. The total contribution is multiplied by a constant factor and is then equally distributed among all, regardless of whether they have contributed initially. Hence, contributors always gain less than free-riders, thus disincentivizing cooperation. In this scenario, arranging a prior commitment or agreement is an essential ingredient in motivating cooperative behaviour, as abundantly observed both in the natural world (Nesse, 2001) and lab experiments (Cherry and McEvoy, 2013).

In (Han et al., 2014), we extend the PGG to examine commitment-based strategies within group interactions. Prior to playing the PGG, commitment-proposing players ask their co-players to commit to contribute to the PGG, paying a personal proposer's cost to establish that agreement. If all of the requested co-players accept the commitment, then the proposers assume everyone will contribute. Those who commit yet later do not contribute must compensate the proposers (Han et al., 2013a). As commitment proposers may encounter non-committers, they require strategies to deal with these individuals. Simplest is to not participate in the creation of the common good. Yet, this avoidance strategy, AVOID, also removes benefits for those wishing to establish the public good, creating a moral dilemma. Alternatively, one can establish boundaries on the common good, so that only those who have truly committed have (better) access, or so that the benefit of non-contributors becomes reduced. This is the RESTRICT strategy. Our results lead to two main conclusions: (i) Both strategies can promote the emergence of cooperation in the one-shot PGG whenever the cost of arranging commitment is justified with respect to the benefit of cooperation, thus generalizing results from pairwise interactions (Han et al., 2013a); (ii) RESTRICT, rather than AVOID, leads to more favourable societal outcomes in terms of contribution level, especially when group size and/or the benefit of the PGG increase, even if the cost of restricting is quite large.

In another approach to commitment-based strategic behaviour in the context of the PGG (Han et al., 2017a), we consider a different set of strategies, envisaging that a restriction measure may not always be possible as it is both costly and takes additional effort to implement. Namely, we consider that before engaging in a group venture agents often secure prior commitments from

other members of the group, and based on the level of participation (i.e. how many group members commit) they can then decide whether or not it is worthwhile to join the group effort (Nesse 2011; Cherry and McEvoy, 2013). This approach is inspired in that many group ventures can be launched only when the majority of the participants do commit to contribute to a common good. Furthermore, while some international agreements require ratification by all parties before entering into force, most (especially global treaties) require some minimum less than the total number of negotiating countries.

We have shown that arranging prior commitments while imposing a minimal participation when interacting in groups can help ensure agents' cooperative behaviour. Namely, our results have shown that if the cost of arranging the commitment is sufficiently small compared to the cost of cooperation, commitment arranging behaviours is frequent, leading thereby to a high level of cooperation in the population. Moreover, an optimal participation level emerges depending both on the dilemma at stake and on the cost of arranging the commitment. Namely, the harsher the common good dilemma is, and the costlier it becomes to arrange the commitment, then the more participants should explicitly commit to the agreement to ensure the success of the joint venture. Furthermore, considering that commitment deals may last for more than one encounter, we evince that longer-lasting commitments require a greater strictness upon fake committers than shorter ones.

In yet another approach to commitment-based strategic behaviour in the context of the PGG (Han et al., 2017b), we consider that agents can delegate the commitment arrangement and participation monitoring processes in the above-described approaches, to a (beneficiary or non-costly) central authority or institution. The institution may itself benefit from improving the level of cooperation in the population or the social welfare (e.g. public transportation arranged by government, international agreements supported by the UN, crowd-sourcing systems) (Nesse 2011; Cherry and McEvoy, 2013). It may also profit directly from this joint activity by requesting a fee from all committed players in order to provide the service. We have shown that this centralised approach to arranging commitments outperforms the described (personalised) commitment strategy. By having a centralised party to help arrange commitments from the group members instead of leaving it to them to have the initiative, it removes the commitment free-riding issue that prevented the personalized approach to achieve full cooperation (Han et al., 2013a, 2017a). As such, it provides a more efficient mechanism for dealing with commitment free-riders, that is those who are not willing to bear the cost of arranging commitments whilst enjoying the overall benefits provided by the paying commitment proposers. We have shown that the participation level plays a crucial role in the decision of whether an agreement should be formed; namely, it needs to be more strict in the centralized system for the agreement to be formed; however, once it is done right, it is much more beneficial in terms of the level of cooperation as well as the attainable level of social welfare.

In short, as commitments have been widely studied in AI and Computer Science, see e.g. (Chopra and Singh, 2009; Chopra et al., 2017), to ensure cooperation in self-organized and distributed multi-agent systems of moral agents, our results have provided important insights into the design of such systems whenever dealing with group interactions. For instance, the key to using the potential of self-organized multi-robot systems is that the robots need to ensure a high level of cooperation amongst themselves, as they may have different skill sets, so as to achieve their tasks successfully. Our group commitment results appear to provide an appropriate method to ensure cooperation: the robots can arrange commitments to warrant that a beneficial coalition of skills is obtained and the task effort is fairly distributed (Sarker et al., 2014). Non-committers may be ostracized from the group and the mission might not be launched if the number of committers is too low. Furthermore, when appropriate (viz. the existence of centralised party to handle the commitment monitoring processes), centralisation of commitments can help deal even better with commitment free-riders. Summing up, ethical fine tunings must be observed when establishing the norms for Public Good Games, whether for humans or non-humans, for otherwise one risks that the inherently desirable common moral

ground may become unfeasible.

## **8. Why is it so hard to say sorry? Commitments bring about sincerity**

When making a mistake, individuals are willing to apologise to secure further cooperation, even if the apology is costly. Similarly, individuals arrange commitments to guarantee that an action such as a cooperative one is in the others' best interest, and thus will be carried out to avoid eventual penalties for commitment failure. Hence, both apology and commitment should proceed side by side in behavioural evolution. In Han et al. (2013b), we studied the relevance of a combination of these two strategies in the context of the IPD. We show that apologising acts are rare in non-committed interactions, especially whenever cooperation is very costly, and that arranging prior commitments can considerably increase the frequency of apologising behaviour. In addition we have shown that, with or without commitments, apology resolves conflicts only if it is sincere, i.e. costly enough. Most interestingly, our model predicts that individuals tend to use a much costlier apology in committed relationships than otherwise, because it helps better identify free-riders, such as fake committers.

Apology is perhaps the most powerful and ubiquitous mechanism for conflict resolution (Abeler et al., 2010; Ohtsubo & Watanabe, 2009), especially among individuals involved in long-term repeated interactions (such as a marriage). An apology can resolve a conflict without having to additionally involve external parties (e.g. teachers, parents, courts), which may cost all sides of the conflict significantly more. Evidence supporting the usefulness of apology abounds, ranging from medical error situations to seller-customer relationships (Abeler et al., 2010). Apology has been implemented in several computerized systems, such as human-computer interaction and online markets, to facilitate users' positive emotions and cooperation (Tzeng, 2004; Utz et al., 2009).

The IPD has been the standard model to investigate conflict resolution and the problem of the evolution of cooperation in repeated interaction settings (Axelrod, 1984; Sigmund, 2010). The IPD game is usually known as a story of tit-for-tat (TFT), which won both Axelrod's tournaments (Axelrod, 1984). TFT cooperates if the opponent cooperated in the previous round, and defects if the opponent defected. But if there can be erroneous moves due to noise (i.e. an intended move is wrongly performed), the performance of TFT declines, because an erroneous defection by one player leads to a sequence of unilateral cooperation and defection. A generous version of TFT, which sometimes cooperates even if the opponent defected (Nowak & Sigmund, 1992), can deal with noise better, yet not thoroughly. For these TFT-like strategies, apology is modelled implicitly as one or more cooperative acts after a wrongful defection.

In (Han et al., 2013b), we describe a model containing strategies that explicitly apologise when making an error between rounds. An apologising act consists in compensating the co-player with an appropriate amount (the higher the more sincere), in order to ensure that this other player cooperates in the next actual round. As such, a population consisting of only apologisers can maintain perfect cooperation. However, other behaviours that exploit this apologetic behaviour could emerge, such as those that accept apology compensation from others but do not apologise when making mistakes (fake apologisers), destroying any benefit of the apology behaviour. Employing EGT (Sigmund, 2010), we show that when the apology occurs in a system where the players first ask for a commitment before engaging in the interaction (Han et al., 2012b, 2012c; Han, 2013), this exploitation can be avoided. Our results lead to these conclusions: (i) Apology alone is insufficient to achieve high levels of cooperation; (ii) Apology supported by prior commitment leads to significantly higher levels of cooperation; (iii) Apology needs to be sincere to function properly, whether in committed relationships or commitment-free ones (which is in accordance with existing experimental studies, e.g. Ohtsubo and Watanabe (2009)); (iv) A much costlier apology tends to be used in committed relationships

than in commitment-free ones, as it can help better identify free-riders such as fake apologisers: “*commitments bring about sincerity*”.

In Artificial Intelligence and Computer Science, apology (Tzeng, 2004; Utz et al., 2009) and commitment (Winikoff, 2007; Wooldridge & Jennings, 1999) have been widely studied, namely how their mechanisms can be formalized, implemented, and used to enhance cooperation in human-computer interactions and online market systems (Tzeng, 2004; Utz et al., 2009), as well as general multi-agent systems (Winikoff, 2007; Wooldridge & Jennings, 1999). Our study provides important insights for the design and deployment of such mechanisms; for instance, what kind of apology should be provided to customers when mistakes are made, and whether apology can be enhanced if complemented with commitments to ensure cooperation, e.g. compensation for customers who suffer wrongdoing.

## **9. Apology and forgiveness evolve to resolve failures in cooperative agreements**

Making agreements on how to behave has been shown to be an evolutionarily viable strategy in one-shot social dilemmas. However, in many situations agreements aim to establish long-term mutually beneficial interactions. Our analytical and numerical results (Martinez-Vaquero et al., 2015, 2017) reveal for the first time under which conditions revenge, apology and forgiveness can evolve, and deal with mistakes within on-going agreements in the context of the IPD. We showed that, when agreement fails, participants prefer to take revenge by defecting in the subsisting encounters. Incorporating costly apology and forgiveness reveals that, even when mistakes are frequent, there exists a sincerity threshold for which mistakes will not lead to the destruction of the agreement, inducing even higher levels of cooperation. In short, even when to err is human, revenge, apology and forgiveness are evolutionarily viable strategies, playing an important role in inducing cooperation in repeated dilemmas.

Using methods from EGT, we provide analytical and numerical insight into the viability of commitment strategies in repeated social interactions, modelled by means of the IPD (Axelrod & Hamilton, 1981). In order to study commitment strategies in the IPD, a number of behavioural complexities need to be addressed. First, agreements may end before the recurring interactions are finished. As such, strategies need to take into account how to behave when the agreement is present and when it is absent, on top of proposing, accepting or rejecting such agreements in the first place. Second, as shown within the context of direct reciprocity (Trivers, 1971), individuals need to deal with mistakes made by an opponent or by themselves, caused for instance by “trembling hands” or “fuzzy minds” (Sigmund, 2010; Nowak, 2006). A decision needs to be made on whether to continue the agreement, or end it collecting the compensation owed from the other’s defection.

As errors might lead to misunderstandings or even the breaking of commitments, individuals may have acquired sophisticated strategies to ensure that mistakes are not repeated or that profitable relationships may continue. Revenge and forgiveness may have evolved exactly to cope with those situations (McCullough, 2008; McCullough et al., 2011). The threat of revenge, through some punishment or withholding of a benefit, may discourage interpersonal harm. Yet, often one cannot distinguish with enough certainty if the other’s behaviour is intentional or just accidental (Han et al., 2011; Fischbacher & Utikal, 2013). In the latter case, forgiveness provides a restorative mechanism that ensures that beneficial relationships can still continue, notwithstanding the initial harm. An essential ingredient for forgiveness, analysed in our work, seems to be (costly) apology (McCullough, 2008), a point emphasised in Smith (2008).

The importance of apology and forgiveness for sustaining long-term relationships has been brought out in different experiments (Abeler et al., 2010; Takaku et al., 2001; Okamoto & Matsumura, 2000; Ohtsubo & Watanabe, 2009). Apology and forgiveness are of interest as they remove the interference of external institutions (which can be quite costly to all parties concerned), in order to ensure cooperation.

Creating agreements and asking others to commit to them provides a basic behavioural mechanism present at all the levels of society, playing a key role in social interactions (Nesse, 2001; Sterelny, 2012; Cherry & McEvoy, 2013). Our work reveals how, when moving to repeated games, the detrimental effect of having a large arrangement cost is moderated, for a subsisting commitment can play its role during several interactions. In these scenarios, the most successful individuals are those who propose commitments (and are willing to pay their cost) and, following the agreement, cooperate unless a mistake occurs. But if the commitment is broken then these individuals take revenge and defect in the remaining interactions, confirming analytically what has been argued in McCullough (2008), and in McCullough et al. (2011). This result is intriguing, since revenge by withholding benefit from the transgressor may lead to a more favourable outcome for cooperative behaviour in the IPD, as opposed to the well-known reciprocal behaviour such as TFT-like strategies. Forgivers only do better when the benefit-to-cost ratio is high enough.

Yet, as mistakes during any (long-term) relationship are practically inevitable, individuals need to decide whether it is worthwhile to end the agreement and collect the compensation when a mistake is made, or whether it is better to forgive the co-player and continue the mutually beneficial agreement. To study this question, the commitment model was extended with an apology-forgiveness mechanism, where apology was defined either as an external or individual parameter in the model. In both cases, we have shown that forgiveness is effective if it takes place after receiving an apology from the co-players. However, to play a promoting role for cooperation, apology needs to be sincere, in other words, the amount proffered in apology has to be high enough (yet not too high), which is also corroborated by recent experimental psychology (McCullough et al., 2014). This extension to the commitment model produces even higher cooperation levels than in the revenge-based outcome. In the opposite case, fake committers that propose or accept a commitment with the intention of taking advantage of the system (defecting and apologising continuously) will dominate the population. In this situation, the introduction of the apology-forgiveness mechanism destroys the increase of the cooperation level that commitments by themselves produce. Thus, there is a lower-limit on how sincere apology needs to be, as below this limit apology and forgiveness even reduce the level of cooperation one could expect from simply taking revenge. It has been shown in previous works that mistakes can induce the outbreak of cheating or intolerant behaviour in society (Martinez-Vaquero & Cuesta, 2013, 2014), and only a strict ethics can prevent them (Martinez-Vaquero & Cuesta, 2014), which in our case would be understood as forgiving only when apology is sincere.

Commitments in repeated interaction settings may take the form of loyalty (Schneider & Weber, 2013; Back & Flache, 2008), which is different from our commitments regarding posterior compensations, for we do not assume a partner choice mechanism. Loyalty commitment is based on the idea that individuals tend to stay with or select partners based on the length of their prior interactions. We go beyond these works by showing that, even without partner choice, commitment can foster cooperation and long-term relationships, especially when accompanied by sincere apology and forgiveness in case mistakes are made.

Ohtsubo's experiment (Ohtsubo & Watanabe, 2009) shows that a costlier apology is better at communicating sincerity, and as a consequence will be more often forgiven. This observation is shown to be valid across cultures (Takaku et al., 2001). In another laboratory experiment (Fischbacher & Utikal, 2013), the authors showed apologies work because they can help reveal the intention behind a wrongdoer's preceding offence. In compliance with this observation, in our model, apology best serves those who intended to cooperate but defect by mistake.

Despite the fact that "to err is human" (Pope, 1711), our research results demonstrate that behaviours like revenge and forgiveness can evolve to cope with mistakes, even when they occur at high rates. Complicating matters is that mistakes are not necessarily intentional, and that even if they are then it might still be worthwhile to continue with a potential mutually

beneficial agreement. Here, a sincerity threshold exists whereby the cost of apologising should exceed the cost of cooperating if the encouragement of cooperation is the actual goal.

### **10. Non-excessive guilt as a means to promote cooperation without external incentives**

Our recent work investigates theoretical models wherein agents express guilt, in order to study the role of guilt in promoting pro-social behaviour (Pereira et al., 2017). Analytical and numerical methods from EGT are employed to identify the conditions under which enhanced cooperation emerges within the context of the Iterated Prisoners Dilemma (IPD). In our work guilt is modelled explicitly as a counter that keeps track of the number of an agent's transgressions, and where a threshold dictates when is it that guilt alleviation, obtained via self-punishment and behaviour change, is required by the agent itself. Such alleviation is costly for the agent experiencing guilt. Hence, it is interesting to identify whether guilt prone strategies can evolve despite this disadvantage, compared to unemotional strategies and in co-presence with these.

We show that when the system consists of agents that resolve their guilt without considering the co-player's own attitude towards guilt alleviation then cooperation does not emerge. In that case, agents expressing no guilt, or having no incentive to alleviate any guilt they might experience, easily dominate the guilt prone agents. When, on the other hand, the guilt prone focal agent entertains that its guilt only needs to be alleviated when a defecting co-player also manifests guilt alleviation, then cooperation may thrive. Indeed, there is no point in feeling guilty and atoning if no one else does. This observation remains consistent for a generalised model as discussed in (Pereira et al., 2017b).

Interest in machine ethics has significantly increased in recent years, wherein a pertinent theme within that context addresses the computational modelling of human emotions like guilt (Pereira and Saptawijaya, 2016; Marsella and Gratch, 2014). But in contrast to our aim and approach, these studies aim to formalise guilt as part of a MAS, such as in virtual agent and cognitive agent systems, for the purpose of regulating social norms or for improving agent decision making and reasoning processes. Our results and approach provide novel insights into the design of such MAS systems; for instance, if agents such as robots or swarms are equipped with the capacity for feeling guilt, even if it might lead to some costly disadvantage, that capacity however drives the system to an overall more cooperative outcome, wherein they are willing to take reparative actions after wrongdoings. This internalised incentive mechanism provides important insights for engineering non-human moral agents (e.g. machines) without the need for constructing costly external incentive devices, such as institutions, to exercise sanctions or to provide rewards.

### **11. Future Work: Emotional and Counterfactual Guilt**

A natural extension of our work on intention recognition, commitment, revenge, apology, and forgiveness involves adding guilt, shame, and confession with surplus apology. We leave shame alone for now as it involves reputation, which we did not address above, so as to concentrate on the more basic model of pairwise interactions, without the intrusion of reputational hearsay. Though both have ostensibly evolved to promote cooperation, we believe that guilt and shame can be treated separately. Guilt is an inward phenomenon that can foster apology, and even spontaneous public confession. Shame is inherently public, and it too may lead to apology and request for forgiveness. Shame, however, hinges on being caught, on failing to deceive, and on a mechanism being in place that lets one fall into disrepute.

From an evolutionary viewpoint, guilt is envisaged as an in-built mechanism that tends to prevent wrongdoing because of internal suffering that pressures an agent to confess when wrongs are enacted, alongside a costlier apology and penance, plus an expectation of forgiveness to alleviate or dispel the guilt-induced suffering.

The hypothesis, consequently, is that the emergence of guilt within a population is evolutionarily advantageous as it represents an extra-costly apology compared to a non-guilty one, enacted as it is in order to decrease the added suffering. We can test this hypothesis by adapting our existing model comprising commitment, revenge, apology, and forgiveness, via piggybacking guilt onto it. To do so, one introduces a present/absent guilt parameter such that, on defection by a guilt-ridden player, not only is thereby increased the probability of apology (confession), but also the player spontaneously pays a costlier apology, as a means to atone internal guilt. On the other hand, the co-player will more readily accept a guilty extra-valued apology, and forgive. In addition, this co-player's attitude, if copied, will contribute to favour his own forgiveness by others in the population, in case his own super-apologetic confession of guilt replaces the standard one of absence of acknowledged guilt. The prediction is that guilt will facilitate and speed-up the emergence of cooperation and overall benefit, in spite of its initial heavier individual cost. One reason behind this prediction is that costs of cooperation are compensated for by the costlier guilt apology paid by others. Another reason is that it is in general more conducive to forgiveness, especially in the border cases where the standard apology is outright insufficient.

We know that guilt is alleviated by private confession, e.g. to a priest or psychotherapist, with cost in prayers or fees, plus the renunciation of past failings. In the context of our research, such ersatz confessions and atonements, precisely by exacting a cost, should render temptation to defect less probable – a preference reversal (Correia, 2016) – with the proceeds appertained to some common good (e.g. in a Public Goods Game or the like, as in through charity).

In summary, future research will attempt to show, by simulation if not analytically, that guilt naturally connects with apology and forgiveness mechanisms because of its emergent evolutionary advantage. It seems not too difficult to incorporate into the present framework, by splitting each strategy into one variant experiencing guilt in case of defection, plus a guiltless one. The population at the start would now contain, instead, an admixture of all of both types, for a given fixed cost and extra cost of guilty apology, plus the usual other parameters, namely a forgiveness threshold. The prediction again is that guilt is evolutionarily advantageous, within a range of the overall parameters defining a starting population composition, via EGT evolution with the usual social imitation of strategies with high payoff success.

This further opens the way to treatment of emotions modelled as strategies, guilt being a widely acknowledged one. It should show that one does not need a specific kind of body (namely an anthropomorphic one) for guilt to serve the role of a moral emotion, useful as it is in population settings where moral cooperation attains good value for all regardless of the means of embodiment.

Finally, counterfactual reasoning (Byrne, 2007; Collins et al., 2004; Pereira & Saptawijaya, 2016a, 2017) could be wielded to prime and tune guilt. Presupposing that the agent can reason counterfactually, e.g. given the by now known sequence of plays by co-players it might reason: “Had I before felt guilty instead, and played according to such guilt, then I would have fared better.” As a consequence, the player would then meta-reflectively (Mendonça, 2016) modify its “feeling level” of guilt for the future.

## **12. Conclusion**

We have summarized the results of our own work, reported and surveyed here, on several fundamental facets concerning the emergence and evolution of cooperation in the collective realm, and have provided references to permit following it up in detail.

We have argued how a multiplicity results from our research employing Evolutionary Game Theory (EGT) modelling and experimentation does profitably lead to important insights into machine ethics, such as the design of moral machines, of multi-agent systems, and of contractual algorithms, plus their potential application in human settings too.

One could further envisage the whole of our above approach as purveying a form of fiction, though recognisably a rather abstract one, yet still adumbrated as per the “Moral Feelings from Rocky Fictional Ground” (John, 2016). Indeed, our abstract mathematical and computational fictional simulations might be construed and stretched to fit a bill whereby such fiction would not necessarily offer theorists of emotion or morality immediate embodied evidence, as in novels, say. In contradistinction, it can possibly offer interesting, challenging and conjectural ideas that might benefit the theorising in these domains.

Notwithstanding, a computer scientist friend bemusedly jokes about our “soap opera” research, what with intention recognition, commitment proposal, defection, guilt, apology, forgiveness, revenge...

### **Acknowledgements**

Profound thanks are due to our co-authors of joint published work herein cited, without which the personal summing up and specific philosophical viewpoint above would not have been possible at all. Alphabetically: Ari Saptawijaya, Francisco C. Santos, Luis Martinez-Vaquero, and Tom Lenaerts. Furthermore, L. M. Pereira acknowledges support from grant FCT/MEC NOVA LINCS PEst UID/CEC/04516/2013. TAH from Teesside URF funding (11200174).

### **References**

- Abeler, Johannes, Juljana Calaki, Kai Andree, and Christoph Basek. “The power of apology.” *Economics Letters* 107, no. 2 (2010): 233-235.
- Axelrod, Robert. *The evolution of cooperation*. Vol. 5145. Basic Books (AZ), 1984.
- Hamilton, William D., and Robert Axelrod. “The evolution of cooperation.” *Science* 211, no. 27 (1981): 1390-1396.
- Back, Istvan, and Andreas Flache. “The adaptive rationality of interpersonal commitment.” *Rationality and Society* 20, no. 1 (2008): 65-83.
- Boden, M. A. “Information, computation, and cognitive science.” *Handbook of the Philosophy of Science* 8 (2008): 749-769.
- Byrne, Ruth M. J. *The rational imagination: How people create alternatives to reality*. MIT press, 2007.
- Charniak, Eugene, and Robert P. Goldman. “A Bayesian model of plan recognition.” *Artificial Intelligence* 64, no. 1 (1993): 53-79.
- Cherry, Todd L., and David M. McEvoy. “Enforcing compliance with environmental agreements in the absence of strong institutions: An experimental analysis.” *Environmental and Resource Economics* (2013): 1-15.
- Chopra, Amit K., and Munindar P. Singh. “Multiagent commitment alignment.” In *Proceedings of The 8th International Conference on Autonomous Agents and Multiagent Systems-Volume 2*, pp. 937-944. International Foundation for Autonomous Agents and Multiagent Systems, 2009.
- Chopra, Amit K., Samuel H. Christie V, and Munindar P. Singh. “Splee: a declarative information-based language for multiagent interaction protocols.” In *Proceedings of the 16th Conference on Autonomous Agents and MultiAgent Systems*, pp. 1054-1063. International Foundation for Autonomous Agents and Multiagent Systems, 2017.
- Collins, John David, Edward Jonathan Hall, and Laurie Ann Paul, eds. *Causation and counterfactuals*. MIT Press, 2004.

- Cohen, Philip R., and Hector J. Levesque. "Intention is choice with commitment." *Artificial intelligence* 42, no. 2-3 (1990): 213-261.
- Correia, Vasco. "Weakness of will and self-control'." *Morality and Emotion* (2016): 35. ISBN: 978-1-138-12130-0, pp. 83-98. London: Routledge.
- Deacon, Terrence W. "The hierarchic logic of emergence: Untangling the interdependence of evolution and self-organization." *Evolution and learning: The Baldwin effect reconsidered* (2003): 273-308.
- Dennett, Daniel Clement. *Sweet dreams: Philosophical obstacles to a science of consciousness*. MiT Press, 2005.
- Fehr, Ernst, and Simon Gächter. "Altruistic punishment in humans." *Nature* 415, no. 6868 (2002): 137-140.
- Fischbacher, Urs, and Verena Utikal. "On the acceptance of apologies." *Games and Economic Behavior* 82 (2013): 592-608.
- Fodor, Jerry A. "Special sciences (or: the disunity of science as a working hypothesis)." *Synthese* 28, no. 2 (1974): 97-115.
- Gaspar, Augusta. "Morality and empathy vs. empathy and morality." *Morality and Emotion* (2016): 62-82. ISBN: 978-1-138-12130-0, pp. 83-98. London: Routledge.
- Han, The Anh. "Intention recognition, commitments and their roles in the evolution of cooperation: from artificial intelligence techniques to evolutionary game theory models." *SAPERE series* 9 (2013).
- Han, The Anh. Emergence of social punishment and cooperation through prior commitments. In proceedings of the 30th AAAI Conference on Artificial Intelligence (AAAI 2016), pp. 2494-2500, Phoenix, Arizona, USA.
- Han, The Anh, and Luis Moniz Pereira. 2013a. "Context-dependent incremental decision making scrutinizing the intentions of others via bayesian network model construction." *Intelligent Decision Technologies* 7, no. 4 (2013): 293-317.
- Han, The Anh, and Luis Moniz Pereira. 2013b. "Intention-based decision making via intention recognition and its applications." In *Human Behavior Recognition Technologies: Intelligent Applications for Monitoring and Security*, pp. 174-211. IGI Global, 2013.
- Han, The Anh, and Luis Moniz Pereira. 2013c. "State-of-the-art of intention recognition and its use in decision making." *AI Communications* 26, no. 2 (2013): 237-246.
- Han The Anh, Pereira, Luis Moniz, and Francisco C. Santos. "Intention recognition promotes the emergence of cooperation." *Adaptive Behavior* 19, no. 4 (2011): 264-279.
- Han, The Anh, Pereira, Luis Moniz, and Francisco C. Santos. 2012a. "Corpus-based intention recognition in cooperation dilemmas." *Artificial Life* 18, no. 4 (2012): 365-383.
- Han, The Anh, Pereira, Luis Moniz, and Francisco C. Santos. 2012b. "Intention recognition, commitment and the evolution of cooperation." In *Evolutionary Computation (CEC), 2012 IEEE Congress on*, pp. 1-8. IEEE, 2012.
- Han, The Anh, Pereira, Luis Moniz, and Francisco C. Santos. 2012c. "The emergence of commitments and cooperation." In *Proceedings of the 11th International Conference on Autonomous Agents and Multiagent Systems-Volume 1*, pp. 559-566. International Foundation for Autonomous Agents and Multiagent Systems, 2012.

- Han, The Anh, Pereira, Luis Moniz, and Tom Lenaerts. "Avoiding or restricting defectors in public goods games?." *Journal of the Royal Society Interface* 12, no. 103 (2015): 20141203.
- Han, The Anh, Pereira, Luis Moniz, Francisco C. Santos, and Tom Lenaerts. 2013a. "Good agreements make good friends." *Scientific reports* 3 (2013).
- Han, The Anh, Pereira, Luis Moniz, Francisco C. Santos, and Tom Lenaerts. 2013b. "Why is it so hard to say sorry? evolution of apology with commitments in the iterated Prisoner's Dilemma." In *Proceedings of the Twenty-Third international joint conference on Artificial Intelligence*, pp. 177-183. AAAI Press, 2013.
- Han, The Anh, Pereira, Luis Moniz, Francisco C. Santos, and Tom Lenaerts. 2015a. "Emergence of cooperation via intention recognition, commitment and apology—a research summary." *AI Communications* 28, no. 4 (2015): 709-715.
- Han, The Anh, Francisco C. Santos, Tom Lenaerts, and Luis Moniz Pereira. 2015b. "Synergy between intention recognition and commitments in cooperation dilemmas." *Scientific reports* 5 (2015).
- Han, The Anh, and Tom Lenaerts. "A synergy of costly punishment and commitment in cooperation dilemmas." *Adaptive Behavior* 24, no. 4 (2016): 237-248.
- Han, The Anh, Pereira, Luis Moniz, and Tom Lenaerts. 2017a. "Evolution of commitment and level of participation in public goods games." *Autonomous Agents and Multi-Agent Systems* 31, no. 3 (2017): 561-583.
- Han, The Anh, Luis Moniz Pereira, Luis A. Martinez-Vaquero, and Tom Lenaerts. 2017b. "Centralized vs. Personalized Commitments and their influence on Cooperation in Group Interactions." In proceedings of the 31st AAAI Conference on Artificial Intelligence (AAAI 2017), pp. 2999-3005, San Francisco, USA.
- Han, The Anh, Ari Saptawijaya, and Luis Moniz Pereira. "Moral reasoning under uncertainty." In *International Conference on Logic for Programming Artificial Intelligence and Reasoning*, pp. 212-227. Springer, Berlin, Heidelberg, 2012.
- Hauser, Marc. *Moral minds: How nature designed our universal sense of right and wrong*. Ecco/HarperCollins Publishers, 2006.
- Hodges, A. "Alan Turing: One of the great philosophers." *Phoenix, London* (1997).
- Hofbauer, Josef, and Karl Sigmund. *Evolutionary games and population dynamics*. Cambridge university press, 1998.
- John, Eileen. "Moral feelings from rocky fictional ground." *Morality and Emotion* (2016): 99. ISBN: 978-1-138-12130-0, pp. 83-98. London: Routledge.
- Levin, J. Functionalism. *The Stanford Encyclopaedia of Philosophy*, (2010), E.N. Zalta (ed.), <http://plato.stanford.edu/archives/sum2010/entries/functionalism/>
- Marsella, Stacy, and Jonathan Gratch. "Computationally modeling human emotion." *Communications of the ACM* 57, no. 12 (2014): 56-67.
- Martinez-Vaquero, Luis A., and José A. Cuesta. "Evolutionary stability and resistance to cheating in an indirect reciprocity model based on reputation." *Physical Review E* 87, no. 5 (2013): 052810.
- Martinez-Vaquero, Luis A., and José A. Cuesta. "Spreading of intolerance under economic stress: Results from a reputation-based model." *Physical Review E* 90, no. 2 (2014): 022805.

- Martinez-Vaquero, Luis A., Han, The Anh, Pereira, Luis Moniz, and Tom Lenaerts. "Apology and forgiveness evolve to resolve failures in cooperative agreements." *Scientific reports* 5 (2015).
- Martinez-Vaquero, Luis A., Han, The Anh, Pereira, Luis Moniz, and Tom Lenaerts. "When agreement-accepting free-riders are a necessary evil for the evolution of cooperation." *Scientific Reports* 7 (2017).
- Mendonça, Dina. "Emotions and akratic feelings." *Morality and Emotion* (2016): 50. ISBN: 978-1-138-12130-0, pp. 83-98. London: Routledge.
- McCullough, Michael. *Beyond revenge: The evolution of the forgiveness instinct*. John Wiley & Sons, 2008.
- McCullough, Michael E., Robert Kurzban, and Benjamin A. Tabak. "Evolved mechanisms for revenge and forgiveness." *Understanding and reducing aggression, violence, and their consequences* (2010): 221-239.
- McCullough, Michael E., Eric J. Pedersen, Benjamin A. Tabak, and Evan C. Carter. "Conciliatory gestures promote forgiveness and reduce anger in humans." *Proceedings of the National Academy of Sciences* 111, no. 30 (2014): 11211-11216.
- McDermott, Drew V. *Mind and mechanism*. MIT Press, 2001.
- Mikhail, J., 2007. Universal moral grammar: Theory, evidence and the future. *Trends in cognitive sciences*, 11(4), pp.143-152.
- Nesse, Randolph M. "Natural selection and the capacity for subjective commitment." *Evolution and the capacity for commitment* (2001): 1-44.
- Nowak, M.A., 2006. Five rules for the evolution of cooperation. *science*, 314(5805), pp.1560-1563.
- Nowak, Martin A., and Karl Sigmund. "Tit for tat in heterogeneous populations." *Nature* 355, no. 6357 (1992): 250-253.
- Ohtsubo, Yohsuke, and Esuka Watanabe. "Do sincere apologies need to be costly? Test of a costly signaling model of apology." *Evolution and Human Behavior* 30, no. 2 (2009): 114-123.
- Okamoto, Kyoko, and Shuichi Matsumura. "The evolution of punishment and apology: an iterated prisoner's dilemma model." *Evolutionary Ecology* 14, no. 8 (2000): 703-720.
- Pereira, Luís Moniz. 2012a. "Turing is among us." *Journal of Logic and Computation* 22, no. 6 (2012): 1257-1277.
- Pereira, Luís Moniz. 2012b. "Evolutionary tolerance." In *Philosophy and Cognitive Science*, pp. 263-287. Springer, Berlin, Heidelberg, 2012.
- Pereira, Luís Moniz and Han The Anh. 2011a. "Intention recognition with evolution prospection and causal bayes networks." In *Computational Intelligence for Engineering Systems*, pp. 1-33. Springer Netherlands, 2011.
- Pereira, Luís Moniz and Han The Anh. 2011b. "Elder care via intention recognition and evolution prospection." In *International Conference on Applications of Declarative Programming and Knowledge Management*, pp. 170-187. Springer, Berlin, Heidelberg, 2009.
- Pereira, Luís Moniz, and Ari Saptawijaya. "Modelling morality with prospective logic." *International Journal of Reasoning-based Intelligent Systems* 1, no. 3-4 (2009): 209-221.

- Pereira, Luís Moniz, and Ari Saptawijaya. "Bridging two realms of machine ethics." In *Programming Machine Ethics*, pp. 159-165. Springer International Publishing, 2016.
- Pereira, Luís Moniz, and Ari Saptawijaya. "Abduction and beyond in logic programming with application to morality." (2015). In: Magnani, L. (Ed.), *IfColog Journal of Logics and their Applications*, Special issue on Abduction, 3(1):37-71. London: College Publications.
- Pereira, Luís Moniz, and Ari Saptawijaya. 2016a. *Programming machine ethics*. Vol. 26. Springer, 2016. ISBN: 978-3-319-29353-0, Berlin: Springer-Verlag.
- Pereira, Luís Moniz, and Ari Saptawijaya. 2016b. "Counterfactuals, logic programming and agent morality." In *Applications of Formal Philosophy*, pp. 25-53. Springer, Cham, 2017., ISBN: 978-3319585055, pp. 25-54, Berlin: Springer.
- Pereira, Luis Moniz, Tom Lenaerts, Martinez-Vaquero, Luis A., and Han The Anh. "Social manifestation of guilt leads to stable cooperation in multi-agent systems." In *Proceedings of the 16th Conference on Autonomous Agents and MultiAgent Systems*, pp. 1422-1430. International Foundation for Autonomous Agents and Multiagent Systems, 2017.
- Pope, Alexander. "An essay on criticism, part II." *Lewis, W. Russel Street, Covent Garden (1711)* (1931).
- Powers, Simon T., Daniel J. Taylor, and Joanna J. Bryson. "Punishment can promote defection in group-structured populations." *Journal of theoretical biology* 311 (2012): 107-116.
- Prinz, Jesse. "Emotions, morality, and identity." *Morality and emotion* (2016): 13-34. ISBN: 978-1-138-12130-0, pp. 83-98. London: Routledge.
- Raihani, Nichola J., and Redouan Bshary. "The reputation of punishers." *Trends in ecology & evolution* 30, no. 2 (2015): 98-103.
- Sadri, Fariba. "Logic-based approaches to intention recognition." In *Handbook of Research on Ambient Intelligence and Smart Environments: Trends and Perspectives*, pp. 346-375. IGI Global, 2011.
- Saptawijaya, Ari and Luis Moniz Pereira. "Logic programming applied to machine ethics." In *Portuguese Conference on Artificial Intelligence*, pp. 414-422. Springer, Cham, 2015. LNCS vol. 9273, ISBN 978-3-319-23485-4. Berlin: Springer-Verlag.
- Saptawijaya, Ari, and Luis Moniz Pereira. 2015a. "The potential of logic programming as a computational tool to model morality." In *A Construction Manual for Robots' Ethical Systems*, pp. 169-210. Springer International Publishing, 2015.
- Saptawijaya, Ari, and Luis Moniz Pereira. 2015b. From Logic Programming to Machine Ethics, in: O. Bendel (Ed.), (2018) *Handbuch Maschinenethik*, Berlin: Springer.
- Sarker, Md Omar Faruque, Torbjørn S. Dahl, Elsa Arcaute, and Kim Christensen. "Local interactions over global broadcasts for improved task allocation in self-organized multi-robot systems." *Robotics and Autonomous Systems* 62, no. 10 (2014): 1453-1462.
- Schneider, Frédéric, and Roberto A. Weber. "Long-term commitment and cooperation." (2013), Tech. Rep., Working Paper Series, University of Zurich, Department of Economics.
- Searle, John R. *The construction of social reality*. Simon and Schuster, 1995.
- Searle, John. *Making the social world: The structure of human civilization*. Oxford University Press, 2010.
- Sigmund, Karl. *The calculus of selfishness*. Princeton University Press, 2010.

- Smith, Nick. *I was wrong: The meanings of apologies*. Cambridge University Press, 2008.
- Sterelny, Kim. *The evolved apprentice*. MIT press, 2012.
- Takaku, Seiji, Bernard Weiner, and Ken-Ichi Ohbuchi. "A cross-cultural examination of the effects of apology and perspective taking on forgiveness." *Journal of Language and Social Psychology* 20, no. 1-2 (2001): 144-166.
- Trivers, Robert L. "The evolution of reciprocal altruism." *The Quarterly review of biology* 46, no. 1 (1971): 35-57.
- Turing, Alan M. "Computing machinery and intelligence." *Mind* 59, no. 236 (1950): 433-460.
- Tzeng, Jeng-Yi. "Toward a more civilized design: studying the effects of computers that apologize." *International Journal of Human-Computer Studies* 61, no. 3 (2004): 319-345.
- Utz, Sonja, Uwe Matzat, and Chris Snijders. "On-line reputation systems: The effects of feedback comments and reactions on building and rebuilding trust in on-line auctions." *International Journal of Electronic Commerce* 13, no. 3 (2009): 95-118.
- Neumann, J. von. "von, Morgenstern O." *Theory of games and economic behavior* 1 (1944). Princeton: Princeton University Press.
- Winikoff, Michael. "Implementing commitment-based interactions." In *Proceedings of the 6th international joint conference on Autonomous agents and multiagent systems*, p. 128. ACM, 2007.
- Wooldridge, Michael, and Nicholas R. Jennings. "The cooperative problem-solving process." *Journal of Logic and computation* 9, no. 4 (1999): 563-592.