



Universiteit
Leiden
The Netherlands

The psychological foundations of reputation-based cooperation

Manrique, H.M.; Zeidler, H.; Roberts, G.; Barclay, P.; Walker, M.; Samu, F.; ... ; Raihani, N.

Citation

Manrique, H. M., Zeidler, H., Roberts, G., Barclay, P., Walker, M., Samu, F., ... Raihani, N. (2021). The psychological foundations of reputation-based cooperation. *Philosophical Transactions Of The Royal Society Of London Series B: Biological Sciences*, 376(1838). doi:10.1098/rstb.2020.0287

Version: Publisher's Version

License: [Creative Commons CC BY 4.0 license](https://creativecommons.org/licenses/by/4.0/)

Downloaded from: <https://hdl.handle.net/1887/3250863>

Note: To cite this publication please use the final published version (if applicable).



Review

Cite this article: Manrique HM, Zeidler H, Roberts G, Barclay P, Walker M, Samu F, Fariña A, Bshary R, Raihani N. 2021 The psychological foundations of reputation-based cooperation. *Phil. Trans. R. Soc. B* **376**: 20200287. <https://doi.org/10.1098/rstb.2020.0287>

Accepted: 2 June 2021

One contribution of 20 to a theme issue 'The language of cooperation: reputation and honest signalling'.

Subject Areas:

behaviour, cognition, ecology, evolution, theoretical biology

Keywords:

cooperation, intention attribution, partner choice, perspective-taking, reputation, social cognition

Author for correspondence:

Héctor M. Manrique
e-mail: manrique@unizar.es

The psychological foundations of reputation-based cooperation

Héctor M. Manrique¹, Henriette Zeidler², Gilbert Roberts³, Pat Barclay⁴, Michael Walker⁵, Flóra Samu⁶, Andrea Fariña⁷, Redouan Bshary⁸ and Nichola Raihani⁹

¹Department of Psychology and Sociology, Universidad de Zaragoza, Campus Universitario de Teruel, Ciudad Escolar, s/n. 44003 Teruel, Spain

²Department of Psychology, School of Life and Health Sciences, Aston University, Birmingham, UK

³Independent Researcher, Newcastle upon Tyne, UK

⁴Psychology, University of Guelph, Guelph, Canada

⁵Department of Zoology and Physical Anthropology, Universidad de Murcia, Murcia, Spain

⁶The Institute for Analytical Sociology, Linköping University, Linköping, Sweden

⁷Social, Economic, and Organizational Psychology, Leiden University, Leiden, The Netherlands

⁸Biology, University of Neuchâtel, Neuchâtel, Switzerland

⁹Experimental Psychology, University College London, London, UK

id HMM, 0000-0002-1943-340X; HZ, 0000-0002-5521-102X; GR, 0000-0001-5954-3243; PB, 0000-0002-7905-9069; MW, 0000-0003-4359-7436; FS, 0000-0002-1215-0984; AF, 0000-0002-7149-6571; RB, 0000-0001-7198-8472; NR, 0000-0003-2339-9889

Humans care about having a positive reputation, which may prompt them to help in scenarios where the return benefits are not obvious. Various game-theoretical models support the hypothesis that concern for reputation may stabilize cooperation beyond kin, pairs or small groups. However, such models are not explicit about the underlying psychological mechanisms that support reputation-based cooperation. These models therefore cannot account for the apparent rarity of reputation-based cooperation in other species. Here, we identify the cognitive mechanisms that may support reputation-based cooperation in the absence of language. We argue that a large working memory enhances the ability to delay gratification, to understand others' mental states (which allows for perspective-taking and attribution of intentions) and to create and follow norms, which are key building blocks for increasingly complex reputation-based cooperation. We review the existing evidence for the appearance of these processes during human ontogeny as well as their presence in non-human apes and other vertebrates. Based on this review, we predict that most non-human species are cognitively constrained to show only simple forms of reputation-based cooperation.

This article is part of the theme issue 'The language of cooperation: reputation and honest signalling'.

1. Introduction

Concern for reputation is a key psychological mechanism for explaining the high levels of cooperation observed in humans. Obtaining a good reputation could lead to downstream benefits via one of two routes: individuals might be more likely to be chosen as a partner (reputation-based partner choice, [1]) or they might be more likely to be rewarded [2] by other individuals ('indirect reciprocity', [3,4]; see [2] for a detailed discussion and comparison). Despite the intensive focus on how cooperation can be theoretically promoted by concern for reputation, these theoretical models have tended to 'black-box' the psychology that underpins decision rules. In this review, we aim to highlight the psychological and cognitive mechanisms that might support reputation-based cooperation in humans. We begin by discussing the ontogeny of reputation-based cooperation in humans, and the cognitive mechanisms that probably underpin the ability to evaluate and manage reputation. We argue that the requirement for these mechanisms might largely preclude the emergence of reputation-based cooperation in other species. We

end by presenting a few examples where reputation-based cooperation in non-human species appears to exist, illustrating how reputation-based cooperation might sometimes be achieved by simpler cognitive means.

2. Reputation-based cooperation in humans and other primates

Reputation-based cooperation relies on two distinct capacities: individuals must be able to evaluate the reputations of others as well as be able to strategically manage their own reputation. The cognition underpinning these two facets of reputation-based cooperation is likely to differ (figure 1). Some evidence suggests that children begin to evaluate others on the basis of their prosociality from a very young age (reviewed in [5] but see [6] for failed replication efforts). Evidence also exists in non-human apes and other primates to suggest that individuals are able to evaluate and choose interaction partners on the basis of observed prosociality ([7–10], but see [11]).

In addition to evaluating others, humans also strategically manage their reputation by behaving more cooperatively when there is a possibility that other individuals will learn about their actions (see meta-analysis by Bradley *et al.* [12]). Observability increases cooperation in many domains, including tax compliance [13]; voter turnout [14]; energy conservation [15]; environmentalism [16]; blood donation [17] and more. Most researchers interpret this increased cooperation as being caused by people's concern for reputation.

However, unlike the ability to evaluate others' reputation, this tendency to strategically manage one's own reputation is not present at all stages of life and instead appears to emerge during development. Although young children (under 2 years old) are known to behave prosocially [18–20], such behaviour appears to stem from an intrinsic motivation to satisfy a partner's needs rather from attempts to strategically manage reputation. Children begin to show a concern for reputation from the age of around five, for example, by refraining from stealing from others if they are observed, or making more generous or fairer donations to recipients when their generosity will be revealed to others [21–23]. Other work has shown that a concern with appearing to be prosocial or fair-minded increases over childhood [24], and that children become especially concerned with self-presentation between the ages of 8 and 11 years [25]. At this age, children are increasingly able to inhibit behaviours that might result in social sanctions [26,27] and attempt to present themselves in a positive light to others. At the same time, children become increasingly sceptical about the intentions of others, particularly when it comes to judging prosocial reputations [28]. Thus, it takes most of childhood for humans to hone their ability to understand how one's actions affect our reputations and to behave strategically so as to curate a positive reputation.

Unlike humans, there is scant evidence that non-human primates attempt to strategically manage their reputation. One recent study found that capuchin monkeys were insensitive to the presence of an observer when deciding whether to share food [29], suggesting that capuchins do not attempt to strategically manage their reputation in this way. Studies in chimpanzees have also yielded null results. For instance, although chimpanzees increase effort in a resource acquisition task when watched by a potential competitor, they do not increase effort when watched by a potential cooperation

partner [30]. In the same task, 4- to 5-year-old children increased their efforts both in the presence of a competitive observer and in the presence of a potential future cooperation partner [30]. Similarly, although 5-year-old children share more and steal less when observed by a peer, chimpanzees are not sensitive to the presence of an observer in the same paradigm [31], see also [22,32].

The findings above suggest that (i) cognitive strategies needed for reputation-based cooperation differ depending on whether we consider evaluation of partners versus managing one's own reputation, and (ii) that managing one's own reputation is likely to depend upon more sophisticated socio-cognitive mechanisms. In what follows, we present four socio-cognitive candidates that may frequently be involved in reputation-based cooperation. Most fundamentally, we propose that an extensive working memory is key to developing the sophisticated forms of reputation management seen in humans. Three additional socio-cognitive abilities derive from working memory that are likely to be involved in reputation-based cooperation. These abilities are: (i) delaying gratification, (ii) understanding others' mental states, and (iii) following and enforcing social norms. We show how these building blocks recruit working memory and how they may impinge upon reputation-based cooperation—as well as distinguishing between the cognition needed for evaluating others' reputations and managing one's own reputation, respectively (figure 1).

3. Cognitive mechanisms supporting reputation-based cooperation

(a) Working memory

Following Fuster [33], we define working memory as a mechanism of temporal integration. Crucially, working memory is not synonymous with short-term memory but rather emphasizes both the reactivation of long-term stored information and the integration of new inputs, both of which are likely to be involved in dynamically evaluating and managing reputation. Working memory can be metaphorically likened to a workstation, a place where information is temporarily held and manipulated. Working memory is engaged whenever sophisticated socio-cognitive calculations are needed, such as appreciating that our own perspectives, beliefs and intentions can differ from those of other individuals, and understanding that an individual's intentions might not be accurately represented by his actions.

The ability to successfully manage one's own reputation might often require individuals to monitor how they appear to others. Such monitoring requires the ability to entertain multiple perspectives simultaneously, which makes burdensome demands on working memory [34]. Successfully managing one's own reputation might also involve mental time travel, which allows individuals to imagine how events might unfold in the future. This ability is also likely to involve working memory [35]. Working memory is also likely to be involved in evaluating the reputations of others, for example, by tracking cooperative behaviours [36] and recalling what happened, with whom and when (episodic memory). The complexity of such tasks can be increased further when individuals compare observed behaviours against normative standards, or against behaviours adopted by other individuals. The all-round use of working memory poses some intriguing questions for

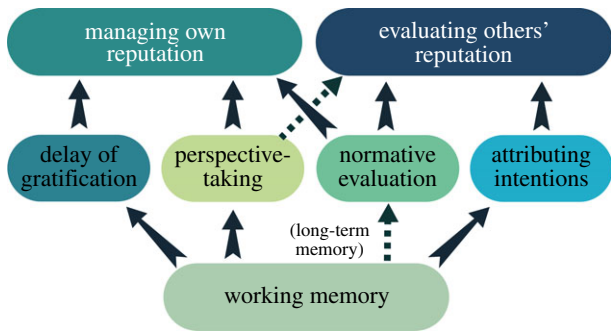


Figure 1. Depiction of how our four socio-cognitive mechanisms are recruited for the managing of one's own reputation as opposed to evaluating third-party reputations. No connecting lines indicate there is no need for the socio-cognitive mechanism in question to be recruited. Arrow continuity expresses the activation of the mechanism is heavily involved in reputation management and/or evaluation. Dotted lines indicate minor involvement. For instance, perspective-taking is key to managing one's own reputation, as we need to see how our acts will appear to a putative observer, yet perspective-taking matters less for evaluating third-party reputations. The opposite is true for attributing intentions. Delay of gratification might be involved in managing one's own reputation as it allows one to resist current temptations to exploit an interaction partner in order to obtain higher future pay-offs associated with curating a good reputation. We expect delay of gratification to be less important for evaluating third-party reputations. Normative understanding is involved in both managing of one's own reputation and evaluating third-party reputations. Working memory is placed in a different level because it enhances the other psychological processes and greatly boosts their efficiency. While working memory is highly involved in delaying gratification, adopting the other's perspective and attributing intentions, its involvement in moral evaluation is lower as norms are stored in long-term memory. (Online version in colour.)

developmental and evolutionary psychology: at what age does children's working memory become capable of maintaining reputation-based cooperative systems? Do great apes have working memory complex enough to sustain reputation-based cooperative systems? By what processes might these abilities have evolved in humans?

Working memory increases linearly between the ages of approximately 7 months and 14 years [37–39]. Meta-analytic evidence [39] suggests that 6-year-olds have a working memory size of three (compared to seven in adults: [40]). Three is the minimum working memory size required to command relative clauses in sentences, which are complex recursive structures like those used to tracking other people's perspectives (e.g. John thinks that Mary knows he is supportive). Given that many reputational acts require such recursion (e.g. John knows that if he does not help Mary now, she will not trust him to reciprocate), it is reasonable to regard three as the minimum working memory size required for constructing complex reputation-based cooperative systems. The extent of working memory involvement in evaluation of others' reputations is likely to vary: evaluations that do not involve recursion (e.g. helping that signals physical ability) may need less working memory than those which do (e.g. helping that signals future intent to cooperate).

Studies directly measuring working memory in great apes are few and have yielded mixed results. Some studies suggest that the working memory capacity of non-human apes is likely to be limited. For instance, in a simplified version of the Wisconsin Card Sorting Test, that involves

sorting cards along three dimensions (shape, colour, number), chimpanzees struggled to form a classificatory criterion or to change it flexibly to match the reinforcement contingencies [41]. Similarly, in a memory task where individuals had to turn over cards one at a time and find matching pairs, chimpanzees made four times more mistakes than humans when tasked with three pairs, which would involve holding three cards in working memory [42]. Nevertheless, other studies have reported remarkable performance in serial ordering tasks administered to chimpanzees, that involved memorizing up to five digits flashed on a screen in ascending order [43], or presenting up to six closed boxes on a platform and having a subject chimpanzee encode and remember those boxes already emptied of food in previous trials to avoid re-opening them again [44].

An alternative approach to assessing working memory capacity involves measuring the extent to which individuals are able to hierarchically classify objects [45–47]. The Langer protocol investigates spontaneous grouping of objects and allows performance to be rated as a function of complexity, ranging from the first-order classifications, where only a single group of objects matching in shape and/or colour is formed (e.g. is set apart from the other objects), to classifications in which more than one group is formed contemporaneously (e.g. rings are grouped together and kept apart from the cubes). Second- (and higher) order classifications are assigned to groups of objects that are perceptually different, yet share the same classificatory criteria. Second- (and higher) order classification impose higher working memory demands on the classificatory rule as well as on the elements to be sorted, as their differing features need to be compared simultaneously and flexibly [45–47]. Chimpanzees attain second-order combinativity around age 5 [48,49] when still they rarely compose more than two sets at a time [47, p. 225]. By contrast, toddlers begin developing three-category classifications around age 3. Three-category classification allows children to hierarchize—such as two subordinate classes within one superordinate class—whereas two-category classification does not [47]. This hierarchization indicates that children develop recursive structures that might help them track other people's perspectives and construct social reputation-based cooperative systems.

Other approaches have inferred working memory size based on the increasing complexity of manufactured stone tools in the fossil record. Making and using simple stone flakes is reported from Late Pliocene Africa 3.4 Ma, where bipedal Australopithecine existed from before 4 Ma. Australopithecines gave rise to the genus *Homo*, perhaps as early as 2.8 Ma, with which they coexisted until after 2 Ma. By 2.5 Ma, there are several Palaeolithic assemblages of sharp conchoidal (i.e. shell-shaped) flakes struck by manual percussion with hard hammer-stones. Conchoidal fracturing requires simultaneously focusing on the core stone, the hammer stone and the percussion angle, which implies a larger working memory than that required for simple flakes [50]. *Homo* predominated by 1.76 Ma, and co-occur in the African archaeological record with flattish stone handaxes. These handaxes often resembled a large almond and were formed by manual percussion with a hard hammer stone that removed small conchoidal flakes in a regular manner (e.g. bifacial stone-tool fashioning), from two surfaces of the handaxe to be. By 0.4–0.3 Ma, handaxes had three-dimensional symmetry, which required their makers to simultaneously remember different perspectives of the core being worked on. To achieve ideal symmetry involves

advanced foresight and the ability to represent mentally the intended final product to exert ongoing corrections on the working substrate. Based on the increasing complexity of stone tools, and the working memory required to make them, a reasonable conjecture is that early *Homo* had a working memory greater than that of Australopithecines, which was in turn similar, if not greater, than that of chimpanzees. Taken together, these various lines of evidence suggest that working memory capacity is likely to be higher in humans than in non-human apes (and specifically chimpanzees).

Although working memory capacity has been relatively understudied in other animals [51], there is some suggestive evidence for correlates of advanced working memory in some species. For example, scrub jays display evidence of episodic-like memory, being able to remember ‘what’, ‘when’ and ‘where’ during food caching events [52], as well as flexibly altering their own caching strategies to avoid being parasitized by others [53]. This example might provide the most compelling evidence for sophisticated working memory in non-primates. As such, if they would benefit from being able to choose partners for cooperative interactions, then they are a good species to test for reputation-based cooperative systems.

(b) Delay of gratification

Any form of costly cooperation based on investments requires the ability to resist the temptation to obtain immediate benefits (e.g. by cheating) in order to pursue a larger benefit in the future. In some cases, this problem may be solved by psychological mechanisms which render cooperative behaviour immediately subjectively rewarding (a phenomenon known as warm glow, [54]). In other cases, individuals may have to effortfully resist an immediately higher-paying option: they must be able to delay gratification.

Although people are systematically present-biased, the human ability to think long term is extraordinary in nature [55,56]. Human consciousness can produce mental simulations of possible futures, allowing decisions to be based on anticipated outcomes [57]. Indeed, a large part of humans’ mental processes seems to be prospective [58], focusing on what ought to be done in the here and now in order to produce positive results in the future [59].

Investing in a prosocial reputation might sometimes require the ability to delay gratification, because the rewards for cooperation come from future (potentially unknown) partners instead of one’s current partner and are therefore inherently more likely to be delayed and less certain to materialize. Several lines of evidence link the ability to delay gratification with cooperative tendency in humans. Focusing on the future makes participants more generous [60], and spurs their willingness to incur personal costs to prevent damaging reputational information from spreading [61]. Children’s ability to delay gratification is positively related to their tendency to share, indicating that the ability to delay gratification might be a prerequisite for children’s sharing and cooperation [62]. Similar patterns have been observed in adults [63,64]; though see [16,65], as well as in blue jays who are prevented from consuming rewards immediately [66]. Children are also better at delaying gratification in cooperative tasks than solo tasks [67]. A direct link between delay of gratification and reputational management has been suggested in 3- and 4-year-old children [68], although other work has shown that people are unable to anticipate the

delayed indirect benefits from their own cooperative investments [69]. To the extent that delay of gratification is involved in reputation-based cooperation, we expect it to be more important in reputation management than in evaluating the reputations of others (figure 1).

In humans, the ability to delay gratification is measured using paradigms such as the ‘marshmallow test’ [70], which measures the willingness to forego a smaller, immediate reward when a larger, delayed reward is promised. Performance on such tasks is variable—and the strategies children use to resist temptation suggest the importance of two different cognitive systems (‘automatic’ versus ‘top-down’) that affect self-control [71,72]. By the age of 6, children become aware that putting the rewards out of sight during the delay interval helps them to withhold and wait longer [73]. By the age of 12, children realize that not only seeing the food influences their performance, but also the way they talk about it—demonstrating the role of metacognition on performance in such settings. Qualitatively similar results have been observed in chimpanzees. In experimental settings, chimpanzees can delay gratification for up to 10 min [74], and seem to use similar strategies to human children to increase performance on these tasks. For example, chimpanzees engage in more play when higher self-restraint is needed in order to gain bigger rewards—suggesting that they are intentionally deploying strategies to increase their performance [75].

The delay-of-gratification test has by now been used on a variety of vertebrate species [76–78] with varying results. Dogs (with their owners) as well as some fishes and large-brained monkeys (macaques and capuchins) are all able to wait for extended periods to obtain larger rewards; cuttlefish have also been reported to wait up to 2 min [79]. By contrast, small monkeys, rats and various birds (pigeons, corvids, parrots) perform poorly in such tasks. Nevertheless, apart from dogs and chimpanzees, individuals of high performing species typically only wait 30–60 s for a larger amount or a preferred food, which offers a stark contrast with the circa 30 min reported in human children [71] in similar tasks—and the potential to delay gratification for much longer periods in adulthood. This reduced delay of gratification in other species may limit their ability to perform reputation-based cooperation.

(c) Theory of mind

Theory of mind is a multifaceted concept that refers to the ability to attribute mental states to oneself and to third parties and encompasses different abilities, which vary in computational complexity. For example, taking another individual’s visual perspective is simpler than attributing intentions, which is in turn simpler than attributing knowledge, which is again simpler than understanding complex perspectives (level 2 perspective-taking) or attributing beliefs. These latter two examples of theory of mind are extremely taxing in terms of computational demands, because they involve entertaining simultaneously alternative, often contradictory, representations of reality (for a more detailed explanation, see [34]).

Here, we introduce two theory of mind abilities that are likely to be involved in reputation management and evaluating the reputation of others: perspective-taking and attribution of intentions. Reputation-based cooperation may be more stable against erosion if bystanders or other third parties can correctly attribute intentions and beliefs to

actors, and if actors can represent how they and their actions are perceived in the eyes of others. For example, an individual may fail to cooperate either because (s)he does not realize that a recipient needs help, or because (s)he currently lacks the resources to help. In other words, individuals with a willingness to help may sometimes behave uncooperatively. If bystanders can correctly identify uncooperative behaviour as a mistake or temporary inability, they can continue a cooperative relationship with those who did not intend to defect. Therefore, the reputation system becomes less prone to errors undermining cooperation.

Errors are particularly problematic in indirect reciprocity models of cooperation. Indirect reciprocity is only stable if agents distinguish between justified defections and unjustified defections (i.e. defecting on defectors versus defecting on cooperators; ‘Kandori’ or ‘standing’ strategies [3,4]. However, such systems are undermined by errors because they can cause two individuals to perceive the same situation differently. Under the Kandori strategy, an actor’s reputation improves if (s)he either helps a partner in good standing or refuses to help a partner in bad standing. Conversely, an actor’s reputation decreases if (s)he fails to help someone in good standing or helps someone in bad standing. Thus, if actors and bystanders evaluate a potential recipient’s reputation differently, bystanders will alter the actor’s reputation score in the opposite direction as the actor (or others) would have expected. Under the Kandori strategy, low frequencies of any type of error may therefore erode cooperation [80]. Perspective-taking (and more broadly theory of mind) are crucial to overcome the limitations of Kandori, as players may acknowledge the possibility of missing information leading to the ‘wrong’ behaviour or the ‘wrong’ interpretation.

By contrast, reputation-based partner choice can function with or without theory of mind. In reputation-based partner choice, actors help others to signal their ability and/or willingness to help [81,82]. Theory of mind is not necessary to signal one’s abilities or to interpret such signals: when people see a good hunter share his kill, they can infer that (s)he is physically skilled enough to catch it (e.g. [83]) without knowing anything of his or her mental state. Hunters need not know anything about the audience’s mental state either—they can learn that certain behaviours are rewarded (e.g. being chosen as a partner) via reinforcement learning. However, theory of mind can greatly aid reputation-based partner choice because it allows for more complex or targeted signals. For example, theory of mind allows audiences to infer a helper’s intentions in order to predict future cooperation and thus allows individuals to signal not just their ability but their willingness to help. Therefore, although simple forms of reputation-based partner choice might be achieved without the advanced socio-cognitive mechanisms we discuss in this paper, we note that reputation-based partner choice can later evolve to become cognitively quite complex, particularly when helpful individuals have an incentive to misrepresent their type to others and when receivers take hidden intentions of partners into consideration when evaluating prosocial acts (see [84] for a detailed discussion).

(i) Perspective-taking

Perspective-taking can be broadly described as the ability to adopt the perspective of others (e.g. visual, informational, emotional). At around 2 years of age, children are able to differentiate what people can or cannot see [85]. However,

it is usually not until 3–4 years of age that children understand that the same item can look different from different perspectives [86]. This ability (level 2 perspective-taking) requires effortful control to suppress the child’s own visual perception, and is often viewed as the precursor to full-blown theory of mind, in which the individual gains the ability to understand others’ knowledge and beliefs.

Perspective-taking is likely to be involved in both reputation management and the evaluation of others’ reputations. Reputation management involves not only behaving in a certain way, but also the ability to shift perspectives to represent how complying or failing to act in this manner will be perceived by others (figure 1). Thus, taking others’ perspectives can make an organism much more effective at reputation management. Similarly, perspective-taking makes an organism better at detecting cheaters: organisms may dishonestly present themselves as cooperative, and it requires cognitive effort for observers to distinguish between genuine versus deceptive cooperators. For example, one individual might normally be a ‘cheater’, but might temporarily act cooperatively when (s)he sees someone (s)he wants to deceive or impress (e.g. a potential mate). Detecting dishonesty involves being able to entertain simultaneously differing views of reality, an ability that can be equated in terms of computational complexity to attributing complex (level 2) visual perspective. Hence, even if perspective-taking is not strictly required to evaluate other’s reputation, managing level 2 visual perspective-taking indicates that organisms have the cognitive potential to entertain simultaneously differing/contrasting views of reality (mine versus yours), and hence the ability of representing simultaneously overt and hidden intentions in other’s actions.

Visual perspective-taking covers a wide spectrum of abilities, from knowing what others can or cannot see (‘level 1’) to understanding that others see something differently as a function of their relative position (level 2) [87,88] and is therefore a good proxy of other mentalizing skills. Level 1 perspective-taking has been extensively investigated in chimpanzees with initially diverging results [89,90]. Karg *et al.* [91] used a variation of the experience projection paradigm [92] where chimpanzees were trained with different pairs of goggles that affected what they could see. When wearing one colour, the apes could see through the goggles but when wearing the other colour, they could not see anything. It could be inferred that chimpanzees are able to shift perspectives if their own experience with the goggles (i.e. seeing versus not seeing) affected their response to human experimenters wearing the goggles. However, in this study, chimpanzees’ gaze-following was not influenced by their own previous experience with the goggles [91]. Subsequent results indicated that chimpanzees may be able to shift perspectives in a competitive context yet correct visual perspective attribution only approached a modest 60% ([91]; but see [93] for more positive findings). Demonstrating level 2 visual perspective-taking in chimpanzees still proves elusive [94].

Outside apes, the basic forms of perspective-taking have currently only been found in large-brained species. For example, rhesus monkeys steal more often from a human competitor whose face is hidden by an opaque barrier than a competitor whose body alone is hidden [95]. Capuchin monkeys can also strategically conceal visual information [95], while macaques have been reported to know what others can or cannot hear [96]. Ravens provide the best evidence for perspective-taking in birds, being able to follow

human gaze direction around obstacles [97] and attributing visual perspectives even to unseen competitors [98]. Most recently, however, there is evidence that cleaner fish *Labroides dimidiatus* females are able to choose foraging sites where their male partners cannot observe them [99]. Altogether, it appears that some other species may have some perspective-taking abilities which can aid reputation-based cooperation, but perhaps not to the same level as humans.

(ii) Attributing intentions

Having a good or bad reputation is not simply the consequence of performing good or bad deeds; the intention behind observed actions matters (although the tendency to take intentions into consideration when forming moral judgements varies across cultures, [100]). Notwithstanding this cross-cultural variability, attributing intentionality is another skill that is key to evaluating third-party reputations (figure 1).

As early as 14 months, infants selectively copy actions performed intentionally, as opposed to those that seem fortuitous [101]. Similarly, Gergely *et al.* [102] showed that 14-month-old children imitate unusual actions (e.g. turning on a light with one's forehead) more often if those actions were voluntary than if the actions were necessary (e.g. the model's hands were full, thus necessitating use of their forehead). Nine to 18-month-old toddlers show more patience towards adults who try but fail to hand them a toy than towards teasing adults (i.e. seem unwilling) [103]. Similarly, 21-month-old children are more willing to help other children who had attempted but failed to hand them a toy in previous interactions, than to those who previously refused to offer the toy [104]. Therefore, it appears children at a very early age can differentiate outcomes from intentions when judging others' behaviour.

Other animals also appear capable of attributing intentions. In one study [105], chimpanzees and orangutans preferentially selected boxes that were deliberately marked as containing rewards, more so than boxes that were accidentally marked by the experimenter. Similar attempts at gauging intention attribution in other non-human primates have met with mixed results: positive in cotton-top tamarins and rhesus macaques [106]; negative in chimpanzees [107], Tonkean macaques and tufted capuchin monkeys [108]. Call *et al.* [109] showed that chimpanzees leave a testing area sooner when confronted with an experimenter who was unwilling to give them food (e.g. a teasing human who took away the food) as opposed to one who was unable to do so. This paradigm has yielded similar results in capuchins and Tonkean macaques [110,111]. Some non-primates also seem able to consider both the intentions and the outcomes of performed actions: grey parrots [112] and even horses [113] behave differently when confronted with an unwilling versus an unable experimenter offering food rewards.

Some intentions are simple and clear, or are even broadcasted, whereas other intentions are hidden—organisms may deliberately hide their intentions in order to trick others. Whereas non-humans may be capable of attributing simple intentions, we think that the ability to represent hidden intentions might be restricted to humans because it might require a full-blown theory of mind, a powerful working memory for simultaneously representing multiple realities or perspectives [34], and possibly even the existence of language for representing knowledge propositionally.

(d) The use of normative rules

The use of norms is a potential key complement to the socio-cognitive abilities discussed in the previous section. Normative/moral understanding is likely to be involved in managing own reputation and in evaluating others' reputations (figure 1). To have a good reputation, individuals must comply with some norms or moral standards and check that their behaviour aligns with those norms. The same goes for judging others' reputations, as individuals must contrast a potential partner's behaviour with the very same normative/moral standards. If humans did not possess an awareness of what the 'right' behaviour is, it would become harder to choose partners based on whether they do the 'right' thing. In indirect reciprocity models, the strong standing strategy makes a clear distinction between what is right and what is wrong, based on the standing of the recipient [3,4]. This can only work if all players converge on a specific norm that defines who is worthy of help, and who is unworthy of help. Thus, indirect reciprocity systems require a species to be able to use norms. By contrast, reputation-based partner choice can function without norms (e.g. if third parties only assess the actor's ability to help). That being said, reputation-based partner choice might also be affected by norms: the same helpful act may be seen as generous if the norm is to help less, or stingy if the norm is to help more [81]. It might be advantageous to compare potential partners to the norm to know whom to choose [114], or to compare oneself to the norm and adjust one's own cooperation up or down accordingly [81,115].

Human infants are born into a world filled with social norms. Throughout infancy, children learn how things are done and not done. By the age of around 2 years, children can follow adults' requests and conform to others' social behaviours [116]. At around the age of 3 years, children can infer norms by observing others acting in a certain way without needing adult directives. At the same time, they also start enforcing norms on others [117]. By around 5 years of age, children reach another milestone of normative development: the spontaneous creation of their own rules [21]. Although cultural norms vary widely in their content and implementation, children all over the world show similar abilities for understanding, following and enforcing socially prescribed behaviours [118]. The ways in which children create and deal with norms suggests a growing understanding that norms are mutual agreements which result in rights and obligations for each individual involved. Interestingly, children's concern about their own reputation (and attempts at actively managing it) seems to trail their normative development [31,119], i.e. children's reputation management develops after their ability to view norms as mutually agreed upon standards for collaborative interactions.

If normative development encompasses the ability to view norms as a set of standards for interactions, then it can only originate in species where collaborative interactions are initiated by joint agency. Given the lack of evidence for shared agency and intentionality in chimpanzees, the existence of a social system based on collective norms and influenced by reputation seems highly unlikely [120,121]. Also, given the sparse evidence for social norms in chimpanzees, it is unsurprising that there is little evidence for norms in other species either. In both vervet monkeys and great tits, there is evidence that migrating individuals may give up previously learned preferences and conform to local arbitrary

preferences [122,123]. If such conformity did represent norm-following, then these species might theoretically be capable of cooperative systems based on social norms. Without such norm-following, the evolution of reputation-based cooperation is less likely or less efficient.

4. Reputation-based cooperation in non-human species

Although cognitive constraints may prevent many non-human species from displaying complex forms of reputation-based cooperation [124], they may have simpler forms that are less cognitively demanding. In social species, individuals often interact in communication networks, where bystanders may eavesdrop on interactions to extract valuable information [125]. Therefore, acting in a communication network has three potential pay-off consequences: the pay-off obtained from the current interaction, the effect of one's own action on the partner's future behaviour towards self and the effect of one's own action on the future behaviour of any bystander that learns about the action. Interactions in a communication network, therefore, allow individuals to identify potentially cooperative or aggressive individuals in their social environment and to adjust their behaviour appropriately. Moreover, the possibility for bystander responsiveness might incentivize individuals to adjust their current behaviour when they are observed, a phenomenon known as 'audience effects' [126]. This concept shares features with reputation management in humans.

While eavesdropping and audience effects are widespread among vertebrates and have even been documented in invertebrates, convincing evidence exists primarily in competitive contexts [125]. By contrast, in species other than our own, there is a paucity of evidence demonstrating that individuals show a concern for gaining a prosocial reputation. Various arguments can be made why signals are likely to be honest in a competitive context [127,128] but less reliable in a cooperative context [81,129–131]. In a competitive context, individual aggressiveness is likely to be correlated with strength, which is based on metastable features like size, muscle mass, agility and experience. Therefore, signals of formability are difficult to fake and more likely to be honest. The honesty of such signals can change the benefits associated with paying attention to them: eavesdropping in order to gain information on a potential partner's formidability is potentially self-serving. In return, strong individuals may benefit from signalling their strength to eavesdropping bystanders, for example, by displaying after a victorious fight, or attacking those lower in the hierarchy after a defeat [132] in order to reduce the likelihood of being the target of future challenges. Strong individuals may even pick a fight that yields a short-term negative pay-off to reduce the likelihood of being challenged by bystanders in the future [128].

Nevertheless, there are a handful of examples from non-human species that are suggestive of reputation-based cooperation. In various species, individuals may temporarily act as a watchman by looking out for predators while the rest of the group forages. While such behaviour has been interpreted as immediately self-serving as it is mostly done by satiated individuals [133], experiments involving dwarf mongooses have shown that playbacks of an individual's watchman calls increases the amount of grooming this

individual receives later in the day [134]. In vervet monkeys, males and females that contribute during territorial disputes receive more grooming by other group members [135]. In Arabian babblers and Siberian jays, males act more aggressively towards predators in the presence of females, which is suggestive of males displaying in the context of female mate choice [136,137]. In all these cases, there is no specific recipient of the initial helpful act, meaning that the source of eventual return benefits is uncertain.

Perhaps, the best studied case is the marine cleaning mutualism involving the cleaner wrasse *Labroides dimidiatus* and its 'client' fish. Cleaners remove ectoparasites from clients, which benefits both partners [138]. However, cleaners prefer to eat client mucus [139], which is detrimental to client health and hence constitutes cheating. As cleaners have about 2000 interactions per day [140], ongoing interactions often take place in the presence of other clients. These bystanders observe the ongoing interaction and invite for inspection if the cleaner behaves cooperatively—but leave if they witness a conflict between cleaner and current client [141], and may swim to another cleaner instead. As a consequence of this client decision rule, cleaners are more cooperative in the presence of bystanders [142,143]. Moreover, cleaners stop adjusting service quality if bystanders stop exerting such partner choice [144,145].

Some features of the cleaner–client interaction structure might facilitate reputation-based cooperation. First, memory requirements are minimal: bystanders need only consider the currently observed interaction to make an immediate decision whether to invite or to avoid inspection. Second, the bystander's decision is self-serving as there is short-term autocorrelation of cleaner service quality; and the clients get immediate feedback on their decisions, which facilitates learning [146]. Cleaners who feed against preference must delay immediate gratification, but the positive or negative feedback of this decision (clients inviting for inspection or swimming away) is almost immediate, which also facilitates learning. Thus, basic reinforcement learning might suffice to achieve reputation-based cooperation in this system.

One obvious distinction between reputation-based cooperation in humans and other animals is that humans use language (see other contributions to this theme issue). Language allows people to flexibly exchange information about other individuals [69]—and can potentially also increase the amount of information that can be exchanged. Language can also help humans to represent (and hence encode) and recall social norms and might also be a prerequisite for expressing more complex aspects of social cognition that are likely to be involved in managing and evaluating reputations. Despite its probable importance, we do not discuss language in this review, because it acts more as a multiplier on other cognitive mechanisms, and we instead focus on other proximate cognitive mechanisms that form the basic building blocks of reputation-based cooperation in humans.

5. Discussion

We have presented four basic psychological building blocks that we consider important facilitators for complex reputation-based cooperation: working memory, delay of gratification, theory of mind and social norms. Working memory allows for parallel processing of diverse information, to properly assess others'

actions and update their reputation scores. Delay of gratification is useful for many types of cooperation, but may be particularly relevant for reputation-based cooperation where the returns come from a future interaction with an observer rather than an immediate reciprocation by one's current partner. Theory of mind makes it easier to properly assess others' actions, and reduces the risk that spreading errors will undermine cooperation. Finally, norms support theory of mind by giving individuals a benchmark of what is right or wrong. The more developed that each of these building blocks is, the more complex the interaction structure can become. We are aware that by picking these four socio-cognitive mechanisms we leave out other processes that might be involved, e.g. long-term memory, yet we think the ones we picked are more critical and better allow for comparison across species.

Reputation-based cooperation based on partner choice might often be less cognitively demanding than that based on indirect reciprocity. On the one hand, reputation-based partner choice might require a better ability to delay gratification (as it might take several acts of investment to outcompete competitors and be chosen by third parties), while indirect reciprocity games are typically set up in such a way that individuals alternate roles as helper and recipient. On the other hand, reputation-based partner choice can exist in cognitively simple forms like 'walk away or reject partner if they seem uncooperative' [114,147]; this does not require high working memory, theory of mind or normative behaviour, though these abilities can make reputation-based partner choice more efficient. By contrast, analyses of indirect reciprocity games have shown that Kandori is the simplest strategy yielding stable cooperation [148], and Kandori requires norms, theory of mind to identify errors and as a consequence more computational power (i.e. working memory). Therefore, the vast majority of animal species may be cognitively constrained from implementing indirect reciprocity [149], and hence be limited to simple forms of reputation-based partner choice. In line with this hypothesis, the few non-human examples of reputation-based cooperation largely fit the concept of reputation-based partner choice, not indirect reciprocity. Most of the examples seem to be about one party gaining information about another, to know whom to cooperate or mate with, or whom to avoid in

fighters—a type of reputation-based partner choice based on eavesdropping [125]. As such, there is a clear evolutionary path for reputation-based partner choice: start with cognitively simple eavesdropping, which then evolves into an active signalling system (see [150] for cues evolving into signals), with more complex abilities arising later in both signallers and receivers in order to perform better within that signalling system.

Future work should further clarify the role of these cognitive mechanisms in reputation-based cooperation in both humans and non-humans. Studies could investigate reputation-based cooperation in humans when these cognitive mechanisms cannot function properly, such as experimental paradigms that increase cognitive load (e.g. [36]), special populations that lack some of these cognitive mechanisms (e.g. [151,152]) or online networks where one cannot use these mechanisms. Non-human studies could artificially grant these abilities to non-humans, for example, by dissociating cooperative investments from ability to delay gratification (cf. [66]). Other studies could use other creative ways of outsourcing cognition to see how they affect reputation-based cooperation. We look forward to seeing further tests of the cognitive building blocks of reputation-based cooperation.

Data accessibility. This article has no additional data.

Authors' contributions. All authors conceived the general idea during the workshop. H.M.M., H.Z. and R.B. wrote a first draft that was then reworked by all authors. The revision of the paper was conducted mainly by H.M.M., R.B., N.R. and P.B.

Competing interests. N.R. is the author of the forthcoming book, *The social instinct: how cooperation shaped the world*.

Funding. H.M.M. wants to thank Prof. Álvaro Arrizabalaga and the MINECO project with reference HAR2017-82483-C3-1-P for financial support. N.R. was supported by a Royal Society University Research Fellowship and by the Leverhulme Trust. P.B. was supported by the Social Science & Humanities Research Council of Canada (SSHRC grant 430287). R.B. was supported by the Swiss Science Foundation (grant no. 310030_192673). A.F. was supported by the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (AdG agreement no. 785635; PI Carsten K.W. De Dreu). F.S. was supported by the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (grant agreement no. 648693, PI: Károly Takács).

References

1. Roberts G. 1998 Competitive altruism: from reciprocity to the handicap principle. *Proc. R. Soc. Lond. B* **265**, 427–431. (doi:10.1098/rspb.1998.0312)
2. Roberts G, Raihani N, Bshary R, Manrique HM, Fariña A, Samu F, Barclay P. 2021 The benefits of being seen to help others: indirect reciprocity and reputation-based partner choice. *Phil. Trans. R. Soc. B* **376**, 20200290. (doi:10.1098/rstb.2020.0290)
3. Kandori M. 1992 Social norms and community enforcement. *Rev. Econ. Stud.* **59**, 63–80. (doi:10.2307/2297925)
4. Ohtsuki H, Iwasa Y. 2007 Global analyses of evolutionary dynamics and exhaustive search for social norms that maintain cooperation by reputation. *J. Theor. Biol.* **24**, 518–531. (doi:10.1016/j.jtbi.2006.08.01)
5. Van de Vondervoort JW, Hamlin JK. 2018 The early emergence of sociomoral evaluation: infants prefer prosocial others. *Curr. Opin. Psychol.* **20**, 77–81. (doi:10.1016/j.copsyc.2017.08.014)
6. Salvadori E, Blazsekova T, Volein A, Karap Z, Tatone D, Mascaro O, Csibra G. 2015 Probing the strength of infants' preference for helpers over hinderers: two replication attempts of Hamlin and Wynn. *PLoS ONE* **10**, e0140570. (doi:10.1371/journal.pone.0140570)
7. Herrmann E, Keupp S, Hare B, Vaish A, Tomasello M. 2013 Direct and indirect reputation formation in nonhuman great apes (*Pan paniscus*, *Pan troglodytes*, *Gorilla gorilla*, *Pongo pygmaeus*) and human children (*Homo sapiens*). *J. Comp. Psychol.* **127**, 63–75. (doi:10.1037/a0028929)
8. Russell YI, Call J, Dunbar RI. 2008 Image scoring in great apes. *Behav. Processes* **78**, 108–111. (doi:10.1016/j.beproc.2007.10.009)
9. Subiaul F *et al.* 2008 Do chimpanzees learn reputation by observation? Evidence from direct and indirect experience with generous and selfish strangers. *Anim. Cogn.* **11**, 611–623. (doi:10.1007/s10071-008-0151-6)
10. Kawai N, Nakagami A, Yasue M, Koda H, Ichinohe N. 2019 Common marmosets (*Callithrix jacchus*) evaluate third-party social interactions of human actors but Japanese monkeys (*Macaca fuscata*) do not. *J. Comp. Psychol.* **133**, 488–495. (doi:10.1037/com0000182)
11. Bueno-Guerra N, Colell M, Call J. 2020 Effects of indirect reputation and type of rearing on food choices in chimpanzees (*Pan troglodytes*). *Behav. Ecol. Sociobiol.* **74**, 1–10. (doi:10.1007/s00265-020-02861-w)

12. Bradley A, Lawrence C, Ferguson E. 2018 Does observability affect prosociality? *Proc. R. Soc. B* **285**, 20180116. (doi:10.1098/rspb.2018.0116)
13. Coricelli G, Joffily M, Montmarquette C, Villeval MC. 2010 Cheating, emotions, and rationality: an experiment on tax evasion. *Exp. Econ.* **13**, 226–247. (doi:10.1007/s10683-010-9237-5)
14. Gerber AA, Green DP, Larimer CW. 2008 Social pressure and voter turnout: evidence from a large-scale field experiment. *Am. Polit. Sci. Rev.* **102**, 33–48. (doi:10.1017/S000305540808009X)
15. Yoeli E, Hoffman M, Rand DG, Nowak MA. 2013 Powering up indirect reciprocity with a large-scale field experiment. *Proc. Natl Acad. Sci. USA* **110**(Suppl. 2), 10 424–10 429. (doi:10.1073/pnas.1301210110)
16. Barclay P, Barker JL. 2020 Greener than thou: people who protect the environment are more cooperative, compete to be environmental, and benefit from reputation. *J. Environ. Psychol.* **72**, 101441. (doi:10.1016/j.jenvp.2020.101441)
17. Lacetera N, Macis M. 2010 Social image concerns and prosocial behavior: field evidence from a nonlinear incentive scheme. *J. Econ. Behav. Organ.* **76**, 225–237. (doi:10.1016/j.jebo.2010.08.007)
18. Dunfield K, Kuhlmeier VA, O'Connell L, Kelley E. 2011 Examining the diversity of prosocial behavior: helping, sharing, and comforting in infancy. *Infancy* **16**, 227–247. (doi:10.1111/j.1532-7078.2010.00041.x)
19. Warneken F, Tomasello M. 2007 Helping and cooperation at 14 months of age. *Infancy* **11**, 271–294. (doi:10.1111/j.1532-7078.2007.tb00227.x)
20. Vaish A, Carpenter M, Tomasello M. 2009 Sympathy through affective perspective taking and its relation to prosocial behavior in toddlers. *Dev. Psychol.* **45**, 534–543. (doi:10.1037/a0014322)
21. Grueneisen S, Tomasello M. 2017 Children coordinate in a recurrent social dilemma by taking turns and along dominance asymmetries. *Dev. Psychol.* **53**, 265–273. (doi:10.1037/dev0000236)
22. Leimgruber KL, Shaw A, Santos LR, Olson KR. 2012 Young children are more generous when others are aware of their actions. *PLoS ONE* **7**, e48292. (doi:10.1371/journal.pone.0048292)
23. McAuliffe K, Blake PR, Warneken F. 2020 Costly fairness in children is influenced by who is watching. *Dev. Psychol.* **56**, 773–782. (doi:10.1037/dev0000888)
24. Shaw A, Montinari N, Piovesan M, Olson KR, Gino F, Norton MI. 2014 Children develop a veil of fairness. *J. Exp. Psychol. Gen.* **143**, 363–375. (doi:10.1037/a0031247)
25. Aloise-Young PA. 1993 The development of self-presentation: self-promotion in 6- to 10-year-old children. *Soc. Cogn.* **11**, 201–222. (doi:10.1521/soco.1993.11.2.201)
26. Apfelbaum EP, Sommers SR, Norton MI. 2008 Seeing race and seeming racist? Evaluating strategic colorblindness in social interaction. *J. Pers. Soc. Psychol.* **95**, 918–932. (doi:10.1037/a0011990)
27. Rutland A, Cameron L, Bennett L, Ferrell J. 2005 Interracial contact and racial constancy: a multi-site study of racial intergroup bias in 3–5 year old Anglo-British children. *J. Appl. Dev. Psychol.* **26**, 699–713. (doi:10.1016/j.appdev.2005.08.005)
28. Heyman G, Barner D, Heumann J, Schenck L. 2014 Children's sensitivity to ulterior motives when evaluating prosocial behavior. *Cogn. Sci.* **38**, 683–700. (doi:10.1111/cogs.12089)
29. Schino G, Boggiani L, Mortelliti A, Pinzaglia M, Addessi E. 2021 Testing the two sides of indirect reciprocity in tufted capuchin monkeys. *Behav. Processes* **182**, 104290. (doi:10.1016/j.beproc.2020.104290)
30. Engelmann JM, Herrmann E, Tomasello M. 2016 The effects of being watched on resource acquisition in chimpanzees and human children. *Anim. Cogn.* **19**, 147–151. (doi:10.1007/s10071-015-0920-y)
31. Engelmann JM, Herrmann E, Tomasello M. 2012 Five-year olds, but not chimpanzees, attempt to manage their reputations. *PLoS ONE* **7**, e48433. (doi:10.1371/journal.pone.0048433)
32. Nettle D, Cronin KA, Bateson M. 2013 Responses of chimpanzees to cues of conspecific observation. *Anim. Behav.* **86**, 595–602. (doi:10.1016/j.anbehav.2013.06.015)
33. Fuster JM. 2001 The prefrontal cortex—an update: time is of the essence. *Neuron* **30**, 319–333. (doi:10.1016/S0896-6273(01)00285-9)
34. Manrique HM, Walker MJ. 2017 *Early evolution of human memory*. Cham, Switzerland: Palgrave Macmillan.
35. Dere D, Zlomuzica A, Dere E. 2019 Fellow travellers in cognitive evolution: co-evolution of working memory and mental time travel? *Neurosci. Biobehav. Rev.* **105**, 94–105. (doi:10.1016/j.neubiorev.2019.07.016)
36. Milinski M, Wedekind C. 1998 Working memory constrains human cooperation in the Prisoner's Dilemma. *Proc. Natl Acad. Sci. USA* **95**, 13 755–13 758. (doi:10.1073/pnas.95.23.13755)
37. Diamond A, Doar B. 1989 The performance of human infants on a measure of frontal cortex function, the delayed-response task. *Dev. Psychobiol.* **22**, 271–294. (doi:10.1002/dev.420220307)
38. Gathercole SE, Pickering B, Ambridge B, Wearing H. 2004 The structure of working memory from 4 to 15 years of age. *Dev. Psychol.* **40**, 177–190. (doi:10.1037/0012-1649.40.2.177)
39. Read DW. 2008 Working memory: a cognitive limit to non-human primate recursive thinking prior to hominid evolution. *Evol. Psychol.* **6**, 147470490800600413. (doi:10.1177/147470490800600413)
40. Miller GA. 1956 The magical number seven, plus or minus two: some limits on our capacity for processing information. *Psychol. Rev.* **63**, 81–97. (doi:10.1037/h0043158)
41. Moriguchi Y, Tanaka M, Itakura S. 2011 Executive function in young children and chimpanzees (*Pan troglodytes*): evidence from a non-verbal dimensional change card sort task. *J. Genet. Psychol.* **172**, 252–265. (doi:10.1080/00221325.2010.534828)
42. Washburn DA, Gullledge JP, James F, Rumbaugh DM. 2007 A species difference in visuospatial working memory: does language link 'what' with 'where'? *Int. J. Comp. Psychol.* **20**, 55–64.
43. Inoue S, Matsuzawa T. 2007 Working memory of numerals in chimpanzees. *Curr. Biol.* **17**, R1004–R1005. (doi:10.1016/j.cub.2007.10.027)
44. Völter CJ, Mundry R, Call J, Seed AM. 2019 Chimpanzees flexibly update working memory contents and show susceptibility to distraction in the self-ordered search task. *Proc. R. Soc. B* **286**, 20190715. (doi:10.1098/rspb.2019.0715)
45. Langer J. 1980 *The origins of logic, six to twelve months*. New York, NY: Academic Press.
46. Langer J. 1986 *The origins of logic: one to two years*. New York, NY: Academic Press.
47. Langer J. 2000 The heterochronic evolution of primate cognitive development. In *Biology, brains and behavior. The evolution of human development* (eds ST Parker, J Langer, ML McKinney), pp. 215–235. Santa Fe, NM: American School of Research Press, 'Advanced Seminar Series', and Oxford, James Currey.
48. Poti P, Langer J, Savage-Rumbaugh S, Brakke KE. 1999 Spontaneous logicomathematical constructions by chimpanzees (*Pan troglodytes P. paniscus*). *Anim. Cogn.* **2**, 147–156. (doi:10.1007/s100710050035)
49. Spinozzi G, Natale F, Langer J, Brakke KE. 1999 Spontaneous class grouping behavior by bonobos (*P. troglodytes*) and common chimpanzees (*P. troglodytes*). *Anim. Cogn.* **2**, 157–170. (doi:10.1007/s100710050036)
50. Read D, Van Der Leeuw S. 2008 Biology is only part of the story. *Phil. Trans. R. Soc. B* **363**, 1959–1968. (doi:10.1098/rstb.2008.0002)
51. Carruthers P. 2013 Evolution of working memory. *Proc. Natl Acad. Sci. USA* **110**(Suppl. 2), 10 371–10 378. (doi:10.1073/pnas.1301195110)
52. Clayton NS, Dickinson A. 1998 Episodic-like memory during cache recovery by scrub jays. *Nature* **395**, 272–274. (doi:10.1038/26216)
53. Correia SP, Dickinson A, Clayton NS. 2007 Western scrub-jays anticipate future needs independently of their current motivational state. *Curr. Biol.* **17**, 856–861. (doi:10.1016/j.cub.2007.03.063)
54. Andreoni J. 1990 Impure altruism and donations to public goods: a theory of warm-glow giving. *Econ. J.* **100**, 464–477. (doi:10.2307/2234133)
55. Roberts WA. 2002 Are animals stuck in time? *Psychol. Bull.* **128**, 473–489. (doi:10.1037/0033-2909.128.3.473)
56. Suddendorf T. 2013 *The gap: the science of what separates us from other animals*. New York, NY: Basic Books.
57. Baumeister RF, Maranges HM, Sjästad H. 2018 Consciousness of the future as a matrix of maybe: pragmatic prospection and the simulation of alternative possibilities. *Psychol. Conscious.: Theory Res. Pract.* **5**, 223–238. (doi:10.31234/osf.io/a3r7h)
58. Seligman MEP, Railton P, Baumeister RF, Sripada C. 2013 Navigating into the future or driven by the past. *Perspect. Psychol. Sci.* **8**, 119–141. (doi:10.1177/1745691612474317)

59. Schacter DL, Addis DR, Buckner RL. 2007 Remembering the past to imagine the future: the prospective brain. *Nat. Rev. Neurosci.* **8**, 657–661. (doi:10.1038/nrn2213)
60. Sjästad H. 2019 Short-sighted greed? Focusing on the future promotes reputation-based generosity. *Judgm. Decis. Mak.* **14**, 199–213. (doi:10.31234/osf.io/fw3gp)
61. Vonasch AJ, Reynolds T, Winegard BM, Baumeister RF. 2018 Death before dishonor: incurring costs to protect moral reputation. *Soc. Psychol. Pers. Sci.* **9**, 604–613. (doi:10.1177/1948550617720271)
62. Sebastián-Enesco C, Warneken F. 2015 The shadow of the future: 5-year-olds, but not 3-year-olds, adjust their sharing in anticipation of reciprocation. *J. Exp. Child Psychol.* **129**, 40–54. (doi:10.1016/j.jecp.2014.08.007)
63. Curry OS, Price ME, Price JG. 2008 Patience is a virtue: cooperative people have lower discount rates. *Pers. Individ. Diff.* **44**, 778–783. (doi:10.1016/j.paid.2007.09.023)
64. Harris AC, Madden GJ. 2002 Delay discounting and performance on the Prisoner's Dilemma game. *Psychol. Rec.* **52**, 429–440. (doi:10.1007/BF03395196)
65. Wu J, Balliet D, Tybur JM, Arai S, Van Lange PA, Yamagishi T. 2017 Life history strategy and human cooperation in economic games. *Evol. Hum. Behav.* **38**, 496–505. (doi:10.1016/j.evolhumbehav.2017.03.002)
66. Stephens DW, McLinn CM, Stevens JR. 2002 Discounting and reciprocity in an iterated Prisoner's Dilemma. *Science* **298**, 2216–2218. (doi:10.1126/science.1078498)
67. Koomen R, Grueneisen S, Herrmann E. 2020 Children delay gratification for cooperative ends. *Psychol. Sci.* **31**, 139–148. (doi:10.1177/0956797619894205)
68. Ma F, Zeng D, Xu F, Compton BJ, Heyman GD. 2020 Delay of gratification as reputation management. *Psychol. Sci.* **31**, 1174–1182. (doi:10.1177/0956797620939940)
69. Wu J, Balliet D, Van Lange PA. 2016 Reputation management: why and how gossip enhances generosity. *Evol. Hum. Behav.* **37**, 193–201. (doi:10.1016/j.evolhumbehav.2015.11.001)
70. Mischel W, Ebbsen EB. 1970 Attention in delay of gratification. *J. Pers. Soc. Psychol.* **16**, 329–337. (doi:10.1037/h0029815)
71. Luerssen A, Gyurak A, Ayduk O, Wendelken C, Bunge SA. 2015 Delay of gratification in childhood linked to cortical interactions with the nucleus accumbens. *Soc. Cogn. Affect. Neurosci.* **10**, 1769–1776. (doi:10.1093/scan/nsv068)
72. Hare TA, Camerer CF, Rangel A. 2009 Self-control in decision-making involves modulation of the vmPFC valuation system. *Science* **324**, 646–648. (doi:10.1126/science.1168450)
73. Mischel HN, Mischel W. 1983 The development of children's knowledge of self-control strategies. *Child Dev.* **54**, 603–619. (doi:10.2307/1130047)
74. Beran MJ, Evans TA. 2006 Maintenance of delay of gratification by four chimpanzees (*Pan troglodytes*): the effects of delayed reward visibility, experimenter presence, and extended delay intervals. *Behav. Processes* **73**, 315–324. (doi:10.1016/j.beproc.2006.07.005)
75. Evans TA, Beran MJ. 2007 Chimpanzees use self-distraction to cope with impulsivity. *Biol. Lett.* **3**, 599–602. (doi:10.1098/rsbl.2007.0399)
76. Miller R, Boeckle M, Jelbert SA, Frohnwieser A, Wascher CA, Clayton NS. 2019 Self-control in crows, parrots and nonhuman primates. *Wiley Interdiscip. Rev. Cogn. Sci.* **10**, e1504. (doi:10.1002/wcs.1504)
77. Susini I, Safryghin A, Hillemann F, Wascher CAF. 2020 Delay of gratification in non-human animals: a review of inter- and intra-specific variation in performance. *BioRxiv*. (doi:10.1101/2020.05.05.078659)
78. Aellen M, Dufour V, Bshary R. 2021 Cleaner fish and other wrasse match primates in their ability to delay gratification. *Anim. Behav.* **176**, 125–143. (doi:10.1016/j.anbehav.2021.04.002)
79. Schnell A.K., Boeckle M, Rivera M, Clayton NS, Hanlon R.T. 2021 Cuttlefish exert self-control in a delay of gratification task. *Proc. R. Soc. B* **288**, 20203161. (doi:10.1098/rspb.2020.3161)
80. Milinski M, Semmann D, Bakker TC, Krambeck HJ. 2001 Cooperation through indirect reciprocity: image scoring or standing strategy? *Proc. R. Soc. Lond. B* **268**, 2495–2501. (doi:10.1098/rspb.2001.1809)
81. Barclay P. 2013 Strategies for cooperation in biological markets, especially for humans. *Evol. Hum. Behav.*, **34**, 164–175. (doi:10.1016/j.evolhumbehav.2013.02.002)
82. Barclay P. 2015 Reputation. In *Handbook of evolutionary psychology*, 2nd edn (ed. D. Buss), pp. 810–828. Hoboken, NJ: J. Wiley & Sons.
83. Smith EA, Bliege Bird R. 2000 Turtle hunting and tombstone opening public generosity as costly signaling. *Evol. Hum. Behav.* **21**, 245–261. (doi:10.1016/S1090-5138(00)00031-3)
84. Raihani N, Power EA. 2021 No good deed goes unpunished: the social costs of prosocial behaviour. *PsyArXiv*. (doi:10.31234/osf.io/ebfgr)
85. Moll H, Tomasello M. 2006 Level 1 perspective-taking at 24 months of age. *Br. J. Dev. Psychol.* **24**, 603–613. (doi:10.1348/026151005X55370)
86. Moll H, Meltzoff AN. 2011 How does it look? Level 2 perspective-taking at 36 months of age. *Child Dev.* **82**, 661–673. (doi:10.1111/j.1467-8624.2010.01571.x)
87. Flavell JH. 1977 The development of knowledge about visual perception. *Nebr. Symp. Motiv.* **25**, 43–76.
88. Flavell JH, Everett BA, Croft K, Flavell ER. 1981 Young children's knowledge about visual perception: further evidence for the Level 1–Level 2 distinction. *Dev. Psychol.* **17**, 99–103. (doi:10.1037/0012-1649.17.1.99)
89. Povinelli DJ, Eddy T. 1996 What young chimpanzees know about seeing. *Monogr. Soc. Res. Child Dev.* **61**, 1–189. (doi:10.2307/1166159)
90. Hare B, Call J, Agnetta B, Tomasello M. 2000 Chimpanzees know what conspecifics do and do not see. *Anim. Behav.* **59**, 771–785. (doi:10.1006/anbe.1999.1377)
91. Karg K, Schmelz M, Call J, Tomasello M. 2015 The goggles experiment: can chimpanzees use self-experience to infer what a competitor can see? *Anim. Behav.* **105**, 211–221. (doi:10.1016/j.anbehav.2015.04.028)
92. Heyes CM. 1998 Theory of mind in nonhuman primates. *Behav. Brain Sci.* **21**, 101–148. (doi:10.1017/S0140525X98000703)
93. Okamoto-Barth S, Call J, Tomasello M. 2007 Great apes' understanding of other individuals' line of sight. *Psychol. Sci.* **18**, 462–468. (doi:10.1111/j.1467-9280.2007.01922.x)
94. Karg K, Schmelz M, Call J, Tomasello M. 2016 Differing views: can chimpanzees do Level 2 perspective-taking? *Anim. Cogn.* **19**, 555–564. (doi:10.1007/s10071-016-0956-7)
95. Flombaum J, Santos L. 2005 Rhesus monkeys can assess the visual perspective of others when competing for food. *Curr. Biol.* **15**, 447–452. (doi:10.1016/j.cub.2004.12.076)
96. Santos LR, Nissen AG, Ferrugia JA. 2006 Rhesus monkeys, *Macaca mulatta*, know what others can and cannot hear. *Anim. Behav.* **71**, 1175–1181. (doi:10.1016/j.anbehav.2005.10.007)
97. Bugnyar T, Stöwe M, Heinrich B. 2004 Ravens, *Corvus corax*, follow gaze direction of humans around obstacles. *Proc. R. Soc. Lond. B* **271**, 1331–1336. (doi:10.1098/rspb.2004.2738)
98. Bugnyar T, Reber SA, Buckner C. 2016 Ravens attribute visual access to unseen competitors. *Nat. Commun.* **7**, 1–6. (doi:10.1038/ncomms10506)
99. McAuliffe K, Drayton LA, Royka A, Aellen M, Santos LR, Bshary R. In press. Do cleaner fish know what others can and cannot see? *Commun. Biol.*
100. Barrett HC *et al.* 2016 Small-scale societies exhibit fundamental variation in the role of intentions in moral judgement. *Proc. Natl Acad. Sci. USA* **113**, 4688–4693. (doi:10.1073/pnas.1522070113)
101. Meltzoff AN. 1995 Understanding the intentions of others: re-enactment of intended acts by 18-month-old children. *Dev. Psychol.* **31**, 838–850. (doi:10.1037/0012-1649.31.5.838)
102. Gergely G, Bekkering H, Király I. 2002 Rational imitation in preverbal infants. *Nature* **415**, 755. (doi:10.1038/415755a)
103. Behne T, Carpenter M, Call J, Tomasello M. 2005 Unwilling versus unable: infants' understanding of intentional action. *Dev. Psychol.* **41**, 328–337. (doi:10.1037/0012-1649.41.2.328)
104. Dunfield KA, Kuhlmeier VA. 2010 Intention-mediated selective helping in infancy. *Psychol. Sci.* **21**, 523–527. (doi:10.1177/0956797610364119)
105. Call J, Tomasello M. 1998 Distinguishing intentional from accidental actions in orangutans (*Pongo pygmaeus*), chimpanzees (*Pan troglodytes*) and human children (*Homo sapiens*). *J. Comp. Psychol.* **112**, 192–206. (doi:10.1037/0735-7036.112.2.192)
106. Wood JN, Glynn DD, Phillips BC, Hauser MD. 2007 The perception of rational, goal-directed action in nonhuman primates. *Science* **317**, 1402–1405. (doi:10.1126/science.1144663)

107. Povinelli DJ, Perilloux HK, Reaux JE, Bierschwale DT. 1998 Young and juvenile chimpanzees' (*Pan troglodytes*) reactions to intentional versus accidental and inadvertent actions. *Behav. Processes* **42**, 205–218. (doi:10.1016/S0376-6357(97)00077-6)
108. Costes-Thiré M, Levé M, Uhrlich P, Pasquarretta C, De Marco A, Thierry B. 2015 Evidence that monkeys (*Macaca tonkeana* and *Sapajus apella*) read moves, but no evidence that they read goals. *J. Comp. Psychol.* **129**, 304–310. (doi:10.1037/a0039294)
109. Call J, Hare B, Carpenter M, Tomasello M. 2004 'Unwilling' versus 'unable': chimpanzees' understanding of human intentional action. *Dev. Sci.* **7**, 488–498. (doi:10.1111/j.1467-7687.2004.00368.x)
110. Canteloup C, Piraux E, Poulin N, Meunier H. 2016 Do Tonkean macaques (*Macaca tonkeana*) perceive what conspecifics do and do not see? *PeerJ* **4**, e1693. (doi:10.7717/peerj.1693)
111. Phillips W, Barnes JL, Mahajan N, Yamaguchi M, Santos LR. 2009 'Unwilling' versus 'unable': capuchin monkeys' (*Cebus apella*) understanding of human intentional action. *Dev. Sci.* **12**, 938–945. (doi:10.1111/j.1467-7687.2009.00840.x)
112. Péron F, Rat-Fischer L, Nagle L, Bovet D. 2010 'Unwilling' versus 'unable': do grey parrots understand human intentional actions?. *Interact. Stud.* **11**, 428–441. (doi:10.1075/is.11.3.06per)
113. Trösch M, Bertin E, Calandrea L, Nowak R, Lansade L. 2020 Unwilling or willing but unable: can horses interpret human actions as goal directed? *Anim. Cogn.* **23**, 1035–1040. (doi:10.1007/s10071-020-01396-x)
114. McNamara JM, Barta Z, Frohman L, Houston AL. 2008 The coevolution of choosiness and cooperation. *Nature* **451**, 189–192. (doi:10.1038/nature06455)
115. Barclay P. 2016 Biological markets and the effects of partner choice on cooperation and friendship. *Curr. Opin. Psychol.* **7**, 33–38. (doi:10.1016/j.copsyc.2015.07.012)
116. Rakoczy H, Schmidt MF. 2013 The early ontogeny of social norms. *Child Dev. Perspect.* **7**, 17–21. (doi:10.1111/cdep.12010)
117. Vaish A, Missana M, Tomasello M. 2011 Three-year-old children intervene in third-party moral transgressions. *Br. J. Dev. Psychol.* **29**, 124–130. (doi:10.1348/026151010X532888)
118. Miller JG. 2007 Cultural psychology of moral development. In *Handbook of cultural psychology* (eds S Kitayama, D Cohen), pp. 477–499. New York, NY: The Guilford Press.
119. Kelsey C, Grossmann T, Vaish A. 2018 Early reputation management: three-year-old children are more generous following exposure to eyes. *Front. Psychol.* **9**, 698. (doi:10.3389/fpsyg.2018.00698)
120. Schmidt MFH, Rakoczy H. 2019 On the uniqueness of human normative attitudes. In *The normative animal? On the anthropological significance of social, moral and linguistic norms* (eds K Bayertz, N Roughley), pp. 121–138. Oxford, UK: Oxford University Press.
121. Tomasello M. 2019 The moral psychology of obligation. *Behav. Brain Sci.* **43**, 1–33. (doi:10.1017/S0140525X19001742)
122. Van de Waal E, Borgeaud C, Whiten A. 2013 Potent social learning and conformity shape a wild primate's foraging decisions. *Science* **340**, 483–485. (doi:10.1126/science.1232769)
123. Aplin LM, Farine DR, Morand-Ferron J, Cockburn A, Thornton A, Sheldon BC. 2015 Experimentally induced innovations lead to persistent culture via conformity in wild birds. *Nature* **518**, 538–541. (doi:10.1038/nature13998)
124. Izuma K. 2012 The social neuroscience of reputation. *Neurosci. Res.* **72**, 283–288. (doi:10.1016/j.neures.2012.01.003)
125. McGregor P. (ed.). 2005 *Animal communication networks*. Cambridge, UK: Cambridge University Press.
126. Matos R, Schlupp I. 2005 Performing in front of an audience: signallers and the social environment. In *Animal communication networks* (ed. P McGregor), pp. 63–83. Cambridge, UK: Cambridge University Press.
127. Arnott G, Elwood RW. 2009 Assessment of fighting ability in animal contests. *Anim. Behav.* **77**, 991–1004. (doi:10.1016/j.anbehav.2009.02.010)
128. Johnstone RA, Bshary R. 2004 Evolution of spite through indirect reciprocity. *Proc. R. Soc. Lond. B* **271**, 1917–1922. (doi:10.1098/rspb.2003.2581)
129. Johnstone RA, Bshary R. 2007 Indirect reciprocity in asymmetric interactions: when apparent altruism facilitates profitable exploitation. *Proc. R. Soc. B* **274**, 3175–3181. (doi:10.1098/rspb.2007.1322)
130. André, J.-B. 2010 The evolution of reciprocity: social types or social incentives? *Am. Nat.* **175**, 197–210. (doi:10.1086/649597)
131. Bebbington K, MacLeod C, Ellison TM, Fay N. 2017 The sky is falling: evidence of a negativity bias in the social transmission of information. *Evol. Hum. Behav.* **38**, 92–101. (doi:10.1016/j.evolhumbehav.2016.07.004)
132. Kazim AJN, Aureli F. 2005 Redirection of aggression: multiparty signalling within a network? In *Animal communication networks* (ed. P McGregor), pp. 191–218. Cambridge, UK: Cambridge University Press.
133. Clutton-Brock TH, O'riain MJ, Brotherton PN, Gaynor D, Kinsky R, Griffin AS, Manser M. 1999 Selfish sentinels in cooperative mammals. *Science* **284**, 1640–1644. (doi:10.1126/science.284.5420.1640)
134. Kern JM, Radford AN. 2018 Experimental evidence for delayed contingent cooperation among wild dwarf mongooses. *Proc. Natl Acad. Sci. USA* **115**, 6255–6260. (doi:10.1073/pnas.1801000115)
135. Arseneau-Robar TJM, Taucher AL, Müller E, van Schaik C, Bshary R, Willems EP. 2016 Female monkeys use both the carrot and the stick to promote male participation in intergroup fights. *Proc. R. Soc. B* **283**, 20161817. (doi:10.1098/rspb.2016.1817)
136. Zahavi A. 1995 Altruism as a handicap: the limitations of kin selection and reciprocity. *J. Avian Biol.* **26**, 1–3. (doi:10.2307/3677205)
137. da Cunha FCR, Fontenelle JCR, Griesser M. 2017 The presence of conspecific females influences male-mobbing behavior. *Behav. Ecol. Sociobiol.* **71**, 52. (doi:10.1007/s00265-017-2267-7)
138. Côté IM. 2000 Evolution and ecology of cleaning symbioses in the sea. *Oceanogr. Mar. Biol.* **38**, 311–355.
139. Grutter AS, Bshary R. 2003 Cleaner wrasse prefer client mucus: support for partner control mechanisms in cleaning interactions. *Proc. R. Soc. Lond. B* **270** (Suppl. 2), S242–S244. (doi:10.1098/rsbl.2003.0077)
140. Grutter AS. 1995 The relationship between cleaning rates and ectoparasites loads in coral reef fishes. *Mar. Ecol. Prog. Ser.* **118**, 51–58. (doi:10.3354/meps118051)
141. Bshary R. 2002 Biting cleaner fish use altruism to deceive image-scoring client reef fish. *Proc. R. Soc. Lond. B* **269**, 2087–2093. (doi:10.1098/rspb.2002.2084)
142. Bshary R, Grutter AS. 2006 Image scoring and cooperation in a cleaner fish mutualism. *Nature* **441**, 975–978. (doi:10.1038/nature04755)
143. Pinto AI, Oates J, Grutter AS, Bshary R. 2011 Cleaner wrasse *Labroides dimidiatus* are more cooperative in the presence of an audience. *Curr. Biol.* **21**, 1140–1144. (doi:10.1016/j.cub.2011.05.021)
144. Triki Z, Wismer S, Levorato E, Bshary R. 2018 A decrease in the abundance and strategic sophistication of cleaner fish after environmental perturbations. *Glob. Change Biol.* **24**, 481–489. (doi:10.1111/gcb.13943)
145. Triki Z, Emery Y, Teles MC, Oliveira RF, Bshary R. 2020 Brain morphology predicts social intelligence in wild cleaner fish. *Nat. Commun.*, **11**, 6423. (doi:10.1038/s41467-020-20130-2)
146. Skinner BF. 1953 *Science and human behavior*. New York, NY: Simon & Schuster.
147. Aktipis CA. 2004 Know when to walk away: contingent movement and the evolution of cooperation. *J. Theor. Biol.* **231**, 249–260. (doi:10.1016/j.jtbi.2004.06.020)
148. Santos FP, Santos FC, Pacheco JM. 2018 Social norm complexity and past reputations in the evolution of cooperation. *Nature* **555**, 242–245. (doi:10.1038/nature25763)
149. Santos FP, Pacheco JM, Santos FC. 2021 The complexity of human cooperation under indirect reciprocity. *Phil. Trans. R. Soc. B* **376**, 20200291. (doi:10.1098/rstb.2020.0291)
150. Biernaskie JM, Perry JC, Grafen A. 2018 A general model of biological signals, from cues to handicaps. *Evol. Lett.* **2**, 201–209. (doi:10.1002/evl3.57)
151. Cage E, Pellicano E, Shah P, Bird G. 2013 Reputation management: evidence for ability but reduced propensity in autism. *Autism Res.* **6**, 433–442. (doi:10.1002/aur.1313)
152. Izuma K, Matsumoto K, Camerer CF, Adolphs R. 2011 Insensitivity to social reputation in autism. *Proc. Natl Acad. Sci. USA* **108**, 17 302–17 307. (doi:10.1073/pnas.1107038108)