

# COGNITIVE AND AFFECTIVE PROCESSES IN CHILDREN'S THIRD-PARTY PUNISHMENT

## Authors

R. Arini, J. Bocarejo Aljure, N. Bueno-Guerra, C. Bayón González, E. Fernández Alba, N. Suárez Fernández, L. Wiggs, B. Kenward

## Abstract

This study investigated how children's punishment affective states change over time, as well as when children begin to prioritise intentions over outcomes in their punishment decisions. Whereas most prior research sampled children from Anglo-America or Northwestern Europe, we tested 5- to 11-year-old children from Colombia and Spain (N = 123). We focused on punishment behaviour in response to ostensibly real moral transgressions, rather than punishment recommendations for hypothetical moral transgressions. We employed moral scenarios involving disloyalty (group-focused moral domain) and unfairness (individual-focused moral domain). Regarding punishment affective states, on average children did not derive much enjoyment from administering punishment, nor did they anticipate that punishment would feel good. Thus, children did not make the same emotional forecasting error adults commonly commit. Regarding the cognitive integration of outcomes and intentions, children began to punish failed intentional transgressions more harshly than accidental transgression, in both disloyalty and unfairness scenarios, much earlier than in previous behavioural studies: around 7 years of age rather than in late adolescence. This could be due to the lower processing demands and higher intention salience of our paradigm. Exploratory analyses revealed that children showed higher concern for disloyalty than unfairness. Punishment of disloyalty remained relatively stable in severity with increasing age, while punishment of unfairness decreased in severity. This suggests that the relative importance of moral concerns for the individual vs. the group may shift because of culture-directed learning processes.

29           **Keywords**

30   Third-party punishment; moral domains; punishment affective states; punishment motives;  
31   outcome-to-intent shift.

32           **Highlights**

- 33       • Children did not derive, or expect to derive, much enjoyment from punishment.
- 34       • Punishment of both unfairness and disloyalty became intention-based around 7 years.
- 35       • Punishment of disloyalty remained stable in severity across children's ages.
- 36       • Punishment of unfairness decreased in severity with children's increasing age.

37

38           **Introduction**

39       Morality consists of a set of norms about how people should or should not behave (Janoff-  
40   Bulman et al., 2009). These norms, in turn, are the product of selective forces driving people to find  
41   solutions to the problems of cooperation that occur in social life (Curry et al., 2019). For  
42   cooperation to be maintained, moral norms need to be enforced. Norm enforcement can take two  
43   main forms: *second-party punishment* (2PP), i.e. punishment of norm transgressors meted out by  
44   the victims; and *third-party punishment* (3PP), i.e. punishment of norm transgressors by unaffected  
45   bystanders who act on behalf of the victims. Whereas second-party punishers correct the behaviour  
46   of transgressors essentially for personal benefits, third-party punishers pay a cost (particularly in  
47   terms of risk of counterretaliation and breakdown of valuable social relationships) for the benefit of  
48   others (Jensen, 2010). 3PP has thus received a great deal of scientific attention given its arguably  
49   altruistic nature (Fehr & Gächter, 2002, but see Raihani & Bshary, 2019). 3PP has been indicated as  
50   a key factor in sustaining the progressive establishment of large-scale cooperative networks in  
51   human societies (Boyd & Richerson, 1992). Indeed, population size and complexity of society have  
52   been shown to predict the level of 3PP (Marlowe et al., 2008).

53 From a developmental perspective, it has been shown that children are willing to enact 3PP  
54 from a very early age (as young as 19 months; Hamlin et al., 2011), in response to a range of norm  
55 transgressions (for a review see Marshall & McAuliffe, 2022). Children engage in 3PP even when it  
56 is costly to do so, whether costs are social (Kenward & Östh, 2015), emotional (Arini et al., 2021;  
57 Yudkin et al., 2020) or economic (Gonzalez-Gadea et al., 2022; McAuliffe et al., 2015; Yang et al.,  
58 2018). Children's 3PP decisions are driven by a variety of motives, concerns and biases: deterrence  
59 of norm transgressors (Arini et al., 2023; Marshall et al., 2021; Twardawski & Hilbig, 2020); justice  
60 restoration (Arini et al., 2023; Riedl et al., 2015); equalisation concerns (Arini et al., 2021; Lee &  
61 Warneken, 2022); intergroup bias (Gummerum et al., 2009; Jordan et al., 2014; Gonzalez-Gadea et  
62 al., 2022); and conformity to a model (Salali et al., 2015; House et al., 2020). However, limited  
63 research has been conducted so far on the emotional experiences of children enacting 3PP, as well  
64 as on the cognitive integration between different types of information into children's 3PP decisions  
65 (reviewed below). This work is thus aimed at shedding light specifically on these two aspects.

### 66 ***Emotional Factors in Punishment***

67 Research about the relation between punishment and emotions has been focused more on the  
68 emotions elicited by moral transgressions (which arguably motivate punishment), rather than on the  
69 emotions elicited by enacting or contemplating punishment. Regarding the former strand of  
70 research, it was found that in adults preference for 2PP was associated with anger towards moral  
71 transgressions, whereas preference for 3PP with disgust (Molho et al., 2017; Tybur et al., 2020).  
72 Additionally, 3PP was predicted by compassion towards the victim (Pfattheicher et al., 2019), and  
73 moral outrage towards the transgressor (Hartsough et al., 2020; Lotz et al., 2011; Ginther et al.,  
74 2022). Thus, it seems that 2PP is consistently elicited by negative emotions, whereas 3PP can be  
75 elicited by both negative and positive emotions. A special case of punishment is represented by  
76 punishment of free riders by the cooperators in the group in a public goods game: since free riding  
77 targets both the self and other group members, punishment combines both 2PP and 3PP. There is

78 evidence that this type of punishment is motivated by anger (Fehr & Gächter, 2002), similarly to  
79 2PP.

80         Developmental studies on the emotions elicited by moral transgressions have focused on the  
81 role of anger, and have demonstrated that the relation between anger and punishment depends on  
82 the interaction between punishers' age and the type of punishment they engage in (2PP vs. 3PP). It  
83 has been shown that violations of both fairness and trustworthiness elicited 2PP, and this  
84 relationship was mediated by anger, from childhood to adulthood (Gummerum et al., 2020; van den  
85 Bos et al., 2012). Violations of fairness also elicited 3PP, but this relationship was mediated by  
86 anger only in adults, not in children or adolescents (Gummerum et al., 2020). Finally, by  
87 experimentally manipulating anger, it was demonstrated that this emotion has a causal role in 2PP  
88 of unfairness in all age groups, whereas in 3PP this occurs only in adults and adolescents, but not in  
89 children (Gummerum et al., 2022).

90         Regarding the emotions elicited by punishment (rather than moral transgressions), studies in  
91 the adult literature indicate that punishment is expected to be experienced as rewarding. Indeed,  
92 adults forecast that punishing uncooperative team members would make them feel better (Carlsmith  
93 et al., 2008). Moreover, people show activation in the striatum (a brain area implicated in reward)  
94 when determining the punishment for those who acted unfairly towards either them or others,  
95 suggesting that they anticipate satisfaction from punishment (De Quervain et al., 2004; Strobel et  
96 al., 2011). By contrast, research about the emotional consequences of punishment has produced  
97 quite mixed results: whereas some studies indicate that enacting punishment induces negative  
98 emotions, others suggest that it can elicit positive emotions under certain conditions. On the one  
99 hand, people who inflicted punishment reported feeling worse than individuals who had not been  
100 given the possibility to punish – an effect mediated by rumination about the transgression suffered  
101 (Carlsmith et al., 2008). On the other hand, seeing the transgressors suffer as a result of punishment  
102 has been shown to have a positive effect on punishers' satisfaction (Eder et al., 2020). However,  
103 seeing the transgressors acknowledge the wrongfulness of their actions had an even stronger effect

104 on punishers' satisfaction (Gollwitzer & Denzler, 2009; Gollwitzer et al., 2011; Aharoni et al.,  
105 2022), since this could be interpreted as a change in moral attitude (Funk et al., 2014). This  
106 evidence thus suggests that adults may perceive punishment as a hedonic experience depending on  
107 how transgressors react to being punished.

108         With respect to the developmental literature, some work suggests that in case of vicarious  
109 2PP and when 3PP is paired with compensation of victims, children may derive enjoyment from  
110 punishment to a certain degree. For example, preschool children have been shown to be willing to  
111 incur costs to watch an agent that had previously mistreated them being punished by someone else  
112 (i.e., vicarious 2PP; Mendes et al., 2018). Although this could be interpreted as evidence that  
113 witnessing punishment is experienced as rewarding, the analysis of children's affective indicators  
114 depicts a more complex picture. Children showed a combination of both positive (i.e., smiles) and  
115 negative emotional expressions (i.e., frowns) while watching the punishment of the antisocial agent  
116 (Mendes et al., 2018), suggesting that they felt both pleasure and distress. Furthermore, when given  
117 the opportunity to themselves respond to transgressions affecting other people, primary school-aged  
118 children reported enjoying enacting 3PP of transgressors, although not as much as compensating  
119 victims (Arini et al., 2023). This may indicate that carrying out both types of behaviours contributes  
120 to children experiencing an overall sense of justice being restored and consequently enjoyment.

121         By contrast, in paradigms in which children could only decide whether to assign 3PP but not  
122 compensation, the emotional consequences of 3PP seem to be consistently negative. More  
123 specifically, children reported experiencing more sadness, less happiness and less excitement when  
124 they engaged in 3PP compared to when they did not (see Supplementary Information in Marshall et  
125 al., 2021). Additionally, children were more likely to report lack of enjoyment when they enacted  
126 real rather than pretend 3PP (Arini et al., 2021, Study 2). This suggests that children's affective  
127 states may be sensitive to the impact of 3PP on transgressors: only when they were really punishing,  
128 but not when they were just pretending to punish, could children have felt responsible for the  
129 suffering of the transgressor. Knowing to be the cause of someone else's suffering may be

130 responsible for children's lack of punishment enjoyment. However, the fact that children enacted  
131 3PP even though they did not find it enjoyable suggests that they may view 3PP as a moral duty to  
132 fulfil for the benefit of others (Arini et al., 2021).

133 Notably, differently to the procedure used with adults by Carlsmith et al. (2008), both  
134 Marshall et al. (2021, Supplementary Information) and Arini et al. (2021, Study 2) asked children to  
135 rate their emotions only after they had already assigned punishment. Therefore, these experimental  
136 paradigms did not rule out the possibility that children decided to carry out 3PP expecting it to be  
137 satisfying, yet they experienced low mood when their expectations were not met (similarly to what  
138 has been found in adults; Carlsmith et al., 2008). To exclude this alternative explanation, in the  
139 present experiment we investigated the temporal changes in 3PP affective states by asking children  
140 to report their affective states before, during and after punishment allocation. We predicted that  
141 neither children's affective states during nor after punishment allocation would be positive, in line  
142 with Marshall et al. (2021, Supplementary Information) and Arini et al. (2021, Study 2). As for  
143 children's affective states before punishment allocation, we made no strong predictions. We  
144 hypothesised that, if children have hedonic expectations about punishment as adults do (Carlsmith  
145 et al., 2008; De Quervain et al., 2004; Strobel et al., 2011), affective states before punishment  
146 allocation would be more positive than those reported during and after punishment allocation. If  
147 instead the thought of carrying out 3PP in isolation consistently evokes negative emotions in  
148 children, affective states before punishment allocation would be no different than those reported  
149 during and after punishment allocation (Table 1, Q1). Moreover, we investigated whether children  
150 are sensitive to the impact of 3PP on transgressors (Arini et al., 2021, Study 2). We hypothesised  
151 that, if children are induced to think about the costs they impose on the transgressors with their 3PP  
152 decisions, they will experience lowering of their affective states due to feeling responsible for the  
153 suffering of the transgressors (Table 1, Q2).

154 ***Integration Between Outcomes and Intentions in Punishment***

155 People's punishment decisions following a moral transgression are affected by the cognitive  
156 integration between different types of information: the *outcome* of the transgressor's action and the  
157 transgressor's *intention* behind such action. Importantly, adults tend to attribute more weight to  
158 intentions over outcomes, across different operationalisations of punishment and study  
159 methodologies (e.g., Barrett et al., 2016; Cushman, 2008; Gummerum & Chu, 2014; Hechler &  
160 Kessler, 2022).

161 Research about the development of children's capability to integrate outcome and intention  
162 information has focused much more on *punishment recommendations* rather than actual punishment  
163 behaviour. This strand of research has made extensive use of vignette tasks, in which children are  
164 presented with hypothetical moral violation scenarios through verbal story-telling, and then asked  
165 whether they consider punishment of norm transgressors an appropriate response. Moral violations  
166 scenarios mostly depict property damage and theft (Baird & Astington, 2004), psychological and  
167 physical harm (Helwig et al., 2001; Nobes et al., 2016; Zelazo et al., 1996), or a combination of the  
168 two (Cushman et al., 2013; Killen et al., 2011; Margoni & Surian, 2017; Martin et al., 2022; Nobes  
169 et al., 2009). Importantly, questions about punishment generally take the form of: "*Should [norm*  
170 *transgressor] get in trouble?*". Since children are not asked whether they themselves would punish  
171 the norm transgressor presented in the vignette, they do not even have to imagine themselves as  
172 hypothetical punishers, but just give an opinion about what would be the right course of action.

173 It has been shown that, when young children are asked to evaluate accidental and failed  
174 intentional transgressions in hypothetical scenarios, the presence of just one negative cue – either  
175 relating to outcomes or intentions – is sufficient for them to recommend punishment. They do not  
176 usually appear to attribute more weight to intentions over outcomes, differently from adults. In fact,  
177 they attribute equal weight to outcomes and intentions (Baird & Astington, 2004; Cushman et al.,  
178 2013; Killen et al., 2011; Margoni & Surian, 2017; Nobes et al., 2016), or more to outcomes over  
179 intentions (Helwig et al., 2001; Martin et al., 2022; Zelazo et al., 1996). It is later on during  
180 development – with the so-called "outcome-to-intent shift" – that children's punishment

181 recommendations tend to become more intention-based. More specifically, condemnation of  
182 accidental transgressions begins to decrease (Cushman et al., 2013), while condemnation of failed  
183 intentional transgressions either remains steady (Cushman et al., 2013) or increases with age  
184 (Martin et al., 2022). Overall, the age of the outcome-to-intent shift for punishment  
185 recommendations varies considerably across studies. A couple of studies found that children as  
186 young as 3 are already able to produce punishment recommendations based on intentions (Nobes et  
187 al., 2009; Van de Vondervoort & Hamlin, 2018). However, the majority of the studies showed that  
188 the outcome-to-intent shift tends to occur in middle childhood, between 5 and 8 years of age (Baird  
189 & Astington, 2004; Cushman et al., 2013; Killen et al., 2011; Martin et al., 2022; Nobes et al.,  
190 2016). It has been proposed that the outcome-to-intent shift may be promoted by both internal  
191 factors, such as the development of theory of mind skills (Killen et al., 2011) and executive  
192 functions (Zelazo et al., 1996), and external factors such as social interactions with adults and peers  
193 (Tomasello et al., 2005). Indeed, understanding others' mental states such as intentions enables  
194 individuals to make predictions about their future behaviours (Young & Tsoi, 2013). This ability, in  
195 turn, may prove crucial to avoid engaging in coordination and negotiation efforts with unreliable  
196 social partners (Grueneisen & Tomasello, 2020, 2022) (see Margoni & Surian, 2016 for a review).

197 More recently, research efforts in developmental psychology have been also directed  
198 towards the investigation of the outcome-to-intent shift in *actual punishment behaviour*. This line of  
199 research has made use of behavioural paradigms (especially economic games such as the ultimatum  
200 game), in which children are required to react to apparently real (rather than hypothetical) moral  
201 violation scenarios, the vast majority of which involve unfair distribution of resources. Most of  
202 these studies focus on 2PP rather than 3PP behaviour. Research on 2PP behaviour has produced  
203 rather mixed results, ranging from no evidence of sensitivity to intentions in early to middle  
204 childhood (Bernhard et al., 2020; Bueno-Guerra et al., 2016; Wittig et al., 2013), to evidence of  
205 sensitivity to intentions already fully developed in primary school-aged children (Jaroslawska et al.,



206 2020; Pelligra et al., 2015; Sutter, 2007) or only emerging during adolescence (Gummerum & Chu,  
207 2014; Güroglu et al., 2009; Güroglu et al., 2011).

208         Regarding instead research on the outcome-to-intent shift in 3PP behaviour, it has been  
209 shown that 4- to 7-year-old children did not differentiate between unequal distributions stemming  
210 from chance or negative intentions (Bernhard et al., 2020), and that both children and adolescents  
211 until 15 years of age consistently based their 3PP responses on outcome information (Gummerum  
212 & Chu, 2014). To date, evidence of the capability to integrate outcomes and intentions in 3PP  
213 behaviour has been found only in adults (Gummerum & Chu, 2014; Hechler & Kessler, 2022),  
214 suggesting that the outcome-to-intent shift in their 3PP behaviour may take place in late  
215 adolescence. However, it has been also shown that, after having witnessed an adult inflicting 3PP  
216 on a norm transgressor, 3- and 4-year-old children were more likely to intervene to reduce the  
217 amount of punishment when the transgressor's misbehaviour was accidental rather than intentional  
218 (Chernyak & Sobel, 2016). This indicates that children may have some degree of sensitivity to  
219 intentions in third-party contexts, even when they are not third-party punishers themselves. Finally,  
220 when we consider partner choice behaviours (i.e., avoiding a norm transgressor could be seen as a  
221 form of indirect punishment), sensitivity to intentions is detectable even in infants: 8-month-olds  
222 preferred to reach for a puppet who was involved in an accidental transgression rather than a failed  
223 intentional transgression (Hamlin, 2013, Study 2).

224         To sum up, the outcome-to-intent shift has been shown to occur, on average, in middle  
225 childhood for punishment recommendations (Baird & Astington, 2004; Cushman et al., 2013;  
226 Killen et al., 2011; Martin et al., 2022; Nobes et al., 2016), and supposedly in late adolescence for  
227 3PP behaviour (Gummerum & Chu, 2014). The developmental lag between expressing intention-  
228 based punishment recommendations and enacting intention-based 3PP behaviour could be due to  
229 the different cognitive demands of different experimental paradigms (vignette tasks vs. behavioural  
230 paradigms; Hilton & Kuhlmeier, 2019). Another, not mutually exclusive explanation is that this  
231 developmental lag is an example of knowledge-behaviour gap (Blake, 2018): children may have

232 beliefs about the right thing to do in response to a moral transgression that they struggle to  
233 implement in practice because of lack of cognitive control skills.

234         Since reducing cognitive demands of tasks has been shown to lower the age at which the  
235 outcome-to-intent occurs in moral judgements (Margoni & Surian, 2020), we took a similar  
236 approach to investigate the outcome-to-intent shift in 3PP behaviour. More specifically, we  
237 developed a behavioural paradigm with arguably lower cognitive demands than the one used by  
238 Gummerum & Chu (2014) to assess whether children can integrate outcome and intention  
239 information into their 3PP behaviour. Whereas in Gummerum & Chu's (2014) paradigm children  
240 had to predict how they would react to a range of possible moral transgressions before observing  
241 them, in our paradigm children were asked to make 3PP decisions after being shown the  
242 transgressions. Moreover, in our paradigm moral scenarios were presented in such a way that  
243 children could infer intentions by observing actors' behaviour and listening to their dialogues as  
244 opposed to having to represent their mental states. By reducing processing demands and increasing  
245 intention salience, we predicted that children would manifest the outcome-to-intent shift in their  
246 3PP behaviour earlier than in late adolescence (Gummerum & Chu, 2014), and potentially within  
247 the same age range commonly observed in punishment recommendations, that is between 5 and 8  
248 years of age (Baird & Astington, 2004; Cushman et al., 2013; Killen et al., 2011; Martin et al.,  
249 2022; Nobes et al., 2016) (Table 1, Q3).

250         Regarding the moral scenarios in our paradigm, we chose unfairness for comparability with  
251 previous literature on the outcome-to-intent shift in 3PP behaviour (Bernhard et al., 2020;  
252 Gummerum & Chu, 2014; Hechler & Kessler, 2022), and disloyalty to assess generalisability of the  
253 findings. This comparison is relevant in light of moral foundations theory, according to which  
254 people's moral concerns pertain to two main domains: an individual-focused domain (including  
255 fairness and harm) aimed at the protection of individuals' rights, and a group-focused domain  
256 (including loyalty, authority and purity) aimed at the formation and maintenance of cohesive social  
257 groups (Graham et al., 2011). Interestingly, it has been found that among adults the role of

258 intentions varies across different types of moral domains: intentions matter more when evaluating  
 259 harm (individual-focused domain) and less when evaluating purity violations (group-focused  
 260 domain) in both US American and British adults (Chakroff et al., 2016; Sweetman & Newman,  
 261 2020a, 2020b; Young & Saxe, 2011; Young & Tsoi, 2013). This finding has been also replicated in  
 262 a large, multi-site study that included even a broad range of small-scale societies, practising  
 263 foraging, pastoralism or horticulture (Barrett et al., 2016). By contrast, nothing is currently known  
 264 about whether the type of moral domain can influence the weight of intentions in children’s 3PP  
 265 behaviour (although speculations have been made by Bernhard et al., 2020). We reasoned that, if  
 266 the pattern of attributing more importance to the role of intentions in individual- over group-focused  
 267 domains is generalisable, children would assign more weight to intentions vs. outcomes for 3PP of  
 268 unfairness than disloyalty. In other words, children would punish failed intentional transgressions  
 269 more severely than accidental transgressions in case of unfairness, but not in case of disloyalty  
 270 (Table 1, Q4).

271 **Table 1. Summary of research questions, associated predictions and whether they were**  
 272 **supported in the present study<sup>1</sup>.**

Topic	Research Question	Prediction	Supported?
<b>Punishment Affective States</b>	Q1: Do children enjoy third-party punishment?	Children do not report positive punishment affective states during and after punishment allocation.	Yes
		No prediction about children’s affective states before punishment allocation.	NA
	Q2: Are children’s punishment affective states influenced by the impact of punishment on transgressors?	Emphasising the impact of punishment on transgressors decreases children’s punishment affective states.	No
<b>Integration Between Outcomes and Intentions in Punishment</b>	Q3: When does the outcome-to-intent shift occur in children’s third-party punishment behaviour?	Children manifest the outcome-to-intent shift in third-party punishment between 5 and 8 years of age.	Yes
	Q4: Do children attribute different weight to intentions vs. outcomes	Children punish failed intentional transgressions more severely than	No

<sup>1</sup> The study was originally additionally intended to examine the effect of the presence or absence of an audience on children’s 3PP severity. Because the manipulation check indicated that the audience manipulation was unsuccessful, we decided to omit discussion of this research question in the main manuscript (see Supplementary Information – sections S1 and S4 for a full description of this variable as used in this study).

---

depending on moral domains when they carry out third-party punishment? accidental transgressions in case of unfairness, but not in case of disloyalty.

---

273

274

275

276

## **Method**

277

### ***Sample***

278

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

The choice of countries for this experiment – Colombia and Spain – was opportunistic but motivated by the desire to counteract sampling bias in developmental psychology (Nielsen et al., 2017; Amir & McAuliffe, 2020), given that the vast majority of studies about punishment mentioned in the Introduction was conducted in Anglo-America or Northwestern Europe. Latin American and Mediterranean societies endorse more collectivist (vs. individualistic) values compared to Anglo-American and Northwestern European societies (Hofstede, 2001), meaning that they place a relatively stronger emphasis on the group (vs. the individual). However, differently from commonly held assumptions, people from Latin American and Mediterranean societies, present a distinctive mixture of independent and interdependent traits in how they relate to others or define themselves. This differentiates them from other collectivist cultures, such as Confucian Asia, where people tend to have more markedly interdependent traits (Krys et al., 2022; Uskul et al., 2023).

We allowed logistical constraints to determine effect sizes; the stopping rule was to collect as much data as possible in the period of time at our disposal. As a result, participants were 123 primary school-aged children, who were tested face-to-face at their schools by the researchers. Of these 123 children, 44 lived in Colombia (*mean age*: 7.7 years; *SD age*: 1.6 years; *age range*: from 5.0 years to 10.8 years; *gender distribution*: 12 girls and 32 boys), and the remaining 79 in Spain (*mean age*: 8.7 years; *SD age*: 1.7 years; *age range*: from 5.3 years to 11.8 years; *gender distribution*: 42 girls and 37 boys). Colombian children were all recruited from the same public

297 school in inner Bogotá and were tested from July 2018 to March 2019. Spanish children were  
298 instead recruited from multiple schools – one mixed public-private school in Oviedo (Asturias), as  
299 well as one public school and two mixed public-private schools in the Madrid region – and tested  
300 from November 2019 to January 2020. Regarding the Colombian sample, all caregivers partially or  
301 fully completed a socio-demographic questionnaire, indicating that they were all of Colombian  
302 nationality, with low-to-middle income and education level (the majority of respondents had a  
303 secondary school qualification). As for the Spanish sample, socio-demographic data was not  
304 systematically collected, but inferred through experimenters’ knowledge of the catchment areas:  
305 caregivers were predominantly of Spanish nationality, with middle-to-high income and education  
306 level. The study was approved by Oxford Brookes University Ethical Review Committee (Study  
307 Number 171101, Children’s Social Judgement in a Computer Game) and received Chair’s approval  
308 by the Universidad de los Andes and Universidad Pontificia Comillas, as well as by the Research  
309 Ethics Committee of the Principality of Asturias.

### 310 ***Materials***

311 We developed a spaceship computer game as a variation of the *MegaAttack* game that had  
312 previously been employed to test British children (Arini et al., 2021, Study 2). The game was  
313 programmed in LÖVE, an open-source game development environment using the LUA  
314 programming language. We then installed the game on various laptop computers that we took to the  
315 test locations to conduct testing sessions in-person. Participants saw on the laptop game bouts that  
316 they were told were being played and commented on live by internet players (but were in fact pre-  
317 recorded). The children’s role was to referee internet players in the *MegaAttack* game, judging  
318 whether they behaved badly or not. In the former case, children could decide whether to assign  
319 punishment to misbehaving internet players and, if so, how much.

320 The game involved a team of two player-controlled spaceships, shooting enemies and  
321 collecting gems. Two in-game tasks subject to potential norm violations were distributing bombs  
322 (one player made the allocation between themselves and their team-member, potentially unfairly), and

323 participating in a cooperative task for the collection of a mega-gem. For the cooperation task to be  
324 successful, both players needed to attach to the mega-gem. If only one of the players attached, with  
325 the other player disloyally ignoring them, the attached player would remain trapped.

326 Each video presenting the moral scenarios via game bouts featured a different pair of player  
327 avatars (different animals inside spaceships) to aid memorisation of their different behaviours, and  
328 was kept short (~1 minute each) with the aim of not excessively taxing children's working memory.  
329 Questions being asked of the children did not require articulated verbal responses. All these  
330 precautions were made to minimise the cognitive demands of our task (Hilton & Kuhlmeier, 2019;  
331 Margoni & Surian, 2020).

332 Regarding the content of the moral scenarios, the **control trial** portrayed cases of moral  
333 norm conformity (i.e., no moral transgressions) in both the fairness and loyalty domains. Therefore,  
334 in the control trial, both outcomes and intentions of the players had the same valence (*positive*  
335 *intention, positive outcome*). Instead, the **test trials** portrayed cases of moral transgressions in either  
336 the fairness or loyalty domains. Moreover, in the test trials, outcomes and intentions of the players  
337 had opposite valences, namely accidental transgressions (*positive intention, negative outcome*) and  
338 failed intentional transgressions (*negative intention, positive outcome*). These two cases are the  
339 most informative to study how the relative weight of intentions and outcome changes with age  
340 (Ingram & Moreno-Romero, 2021). In addition to that, the fact that each video in the test trials  
341 contained only one negative cue, either relating to outcomes or intentions, ruled out the potential  
342 inconvenience that children could merely anchor their 3PP decisions to the first negative cue  
343 appearing in the scenarios (Nelson, 1980).

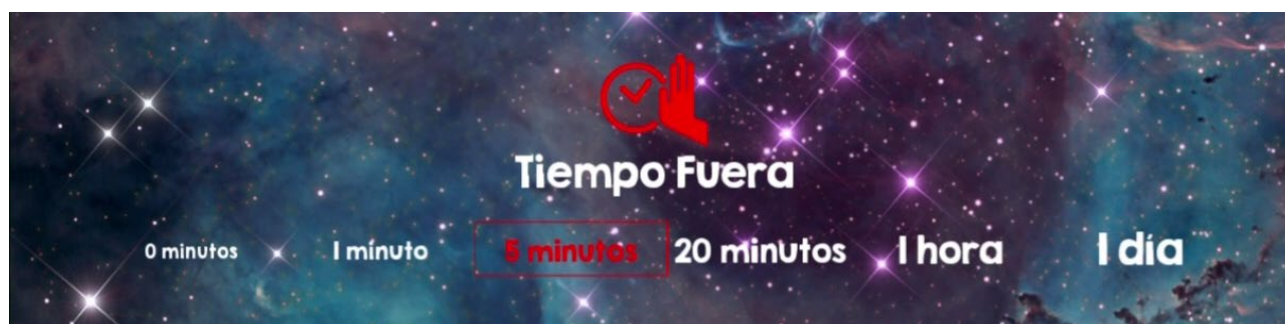
#### 344 ***Design***

345 We adopted a mixed design in which the factors were: *Moral domain* (2 within-subject  
346 levels: fairness domain; loyalty domain); *Intentionality* (3 within-subject levels: failed intentional  
347 transgression; accidental transgression; no moral transgression); *Question time* (3 within-subject

348 levels: before; during; after); *Question focus* (2 between-subject levels: focus on 3PP impact; no  
349 focus on 3PP impact).

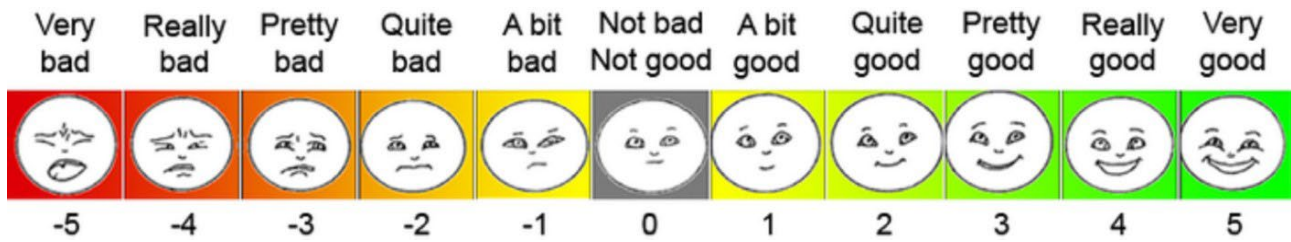
350 All children participating in the experiment were presented first with the control trial  
351 (portraying no moral transgressions), followed by the four test trials (portraying moral  
352 transgressions) in counterbalanced order. Order with respect to failed intentional/accidental  
353 transgression was ABBA or BAAB, and with respect to disloyalty/unfairness transgression was  
354 ABAB or BABA (see Supplementary Information – Table S1). Notably, by consistently showing  
355 the control trial at the beginning of the refereeing sessions, we ensured that participants were always  
356 exposed to the same reference point against which to compare subsequent trials (see e.g.  
357 Twardawski & Hilbig, 2020 and Arini et al., 2021 for similar design choices). We reasoned that this  
358 approach would aid children’s understanding of what was expected of them as referees (i.e.,  
359 allocating punishment only in case of moral transgressions).

360 The dependent variables measured were: *Punishment severity* (6 ordinal levels ranging from  
361 1, “no punishment”, to 6, “1 day-ban”, Figure 1); *Punishment affective states* (11 ordinal levels  
362 from -5, “very bad”, to +5, “very good”, Figure 2). We also measured for use as a control covariate  
363 *Judgement of transgression severity* (6 ordinal levels, ranging from -5, “very bad” to 0, “not bad not  
364 good”, Figure 2), given that substantial variance in punishment severity and affective states has  
365 been shown to be explained by this variable (Arini et al., 2021, 2023).



366  
367 **Figure 1. Scale used to measure punishment severity.** Punishment was a time-out from the game;  
368 children were asked to decide the length of the time-out (“How long do you want the time out to  
369 be?”) among the following options: 0 minutes; 1 minute; 5 minutes; 20 minutes; 1 hour; 1 day.

370  
371



372

373

374

375

376

377

378

379

380

381

381

**Figure 2. Scale used to measure both judgments of transgression severity and punishment affective states.** When children thought that a player had misbehaved in the game, they were asked to judge the severity of the player’s transgression (“*How bad do you think that was?*”). Their options ranged from “very bad” to “not bad, not good” (first 6 points of the scale). Additionally, children were asked to rate their punishment feelings at three time points: before (“*How do you think it will feel to do that?*”), during (“*How did it feel to do that?*”) and after punishment allocation (“*How did it make you feel?*”). The options they were given to rate their feelings ranged from “very bad” to “very good” (all 11 points of the scale).

382

### **Procedure**

383

384

385

386

387

388

389

390

391

392

393

394

The experimental phases were playing familiarisation, refereeing introduction, refereeing sessions, and manipulation check questions. In the **playing familiarisation** the experimenter and the child played together as a team, and the experimenter illustrated to the child the moral norms applied to the *MegaAttack* game. Team-mates were expected to equally divide some bombs among themselves (i.e., fairness norm), and to offer each other help in a cooperative task for the collection of a mega-gem (i.e., loyalty norm). In the **refereeing introduction** the child was told that they would switch from the role of player to that of referee in the *MegaAttack* game. In this new role, the child would have to judge the behaviour of some internet players, with the possibility to give them a time-out from the game when a moral transgression had occurred. We chose this form of punishment for its ecological validity: real computer games have implemented similar systems that give players the possibility to punish misbehaving players by temporarily or permanently banning them from the game (Kou et al., 2017).

395

396

397

398

399

In the actual **refereeing sessions** of the experiment (consisting of one control trial and four test trials), the refereeing child watched five purportedly live game bouts which had been actually pre-recorded. During these game bouts, the child could hear dialogues between the internet players describing their own intentions; gender of these voice-overs was matched with that of the child being tested (video rendition of the refereeing sessions with dialogues in Spanish and English



400 translation available at the Open Science Framework: <https://osf.io/c9w2a/>). In the **control trial** the  
401 child watched one game bout in which no moral norms were violated by the two internet players.  
402 Since the players were both loyal and fair to each other, the refereeing child was expected to  
403 conclude that no misbehaviours had occurred.

404 In the **test trials** the child watched four game bouts, representing a combination of norm  
405 transgressions varying in terms of moral domain and intentionality (Supplementary Information –  
406 section 1.6). Regarding the norm transgressions being shown in the videos, they could be either  
407 accidental transgressions or failed intentional transgressions, related either to the fairness or loyalty  
408 domain. Accidental transgressions were characterised by players having positive intentions,  
409 followed by negative outcomes. Conversely, failed intentional transgressions were characterised by  
410 negative intentions followed by positive outcomes. More specifically, in **accidental unfairness**,  
411 one player intended to split the bombs equally with the team-member (5 bombs each, out of 10) but,  
412 by mistake, ended up with more bombs (7/10) than the equal share (Figure 3A). In **failed**  
413 **intentional unfairness**, one player intended to take for themselves more bombs (7/10) than the  
414 equal share, but inadvertently ended up allocating equal numbers of bombs (5/10) to themselves and  
415 the team-member (Figure 3B). In **accidental disloyalty**, one player intended to cooperate with the  
416 team-mate in the mega-gem collection but, due to a mistake, failed to free the trapped team-mate  
417 from the mega-gem (Figure 3C). In **failed intentional disloyalty**, one player intended to leave the  
418 team-mate trapped in the mega-gem, but inadvertently set them free (Figure 3D).

419 After having seen each of the five game bouts, the child had to answer for each of the two  
420 players in turn: “*Did this player behave badly?*”. If a misbehaviour was identified, the child had to  
421 express their judgement of transgression severity of the norm transgression (“*How bad do you think*  
422 *that was?*”; answers provided by using the scale in Figure 2). The child was then asked to establish  
423 the punishment severity for the norm transgressor, operationalised as a time-out from the game  
424 (“*How long do you want the time out to be?*”; answers provided by using the scale in Figure 1).  
425 Additionally, children were asked to rate their punishment affective states at three time points

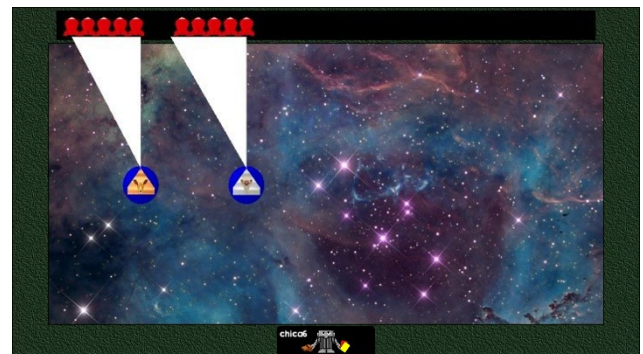
426 (answers provided by using the scale in Figure 2): 1) before their first 3PP decision by forecasting  
427 how punishment would feel (time point: before; question: “*How do you think it will feel to do*  
428 *that?*”); 2) once they had punished for the first time (time point: during; question: “*How did it feel*  
429 *to do that?*”); 3) after the last 3PP decision (time point: after; question: “*How did it make you*  
430 *feel?*”). The sentence preceding the question about punishment affective states was manipulated in  
431 order to either highlight, or not, the cost suffered by the transgressors due to children’s  
432 administration of punishment (between-subjects framing manipulation: focus on 3PP impact vs. no  
433 focus on 3PP impact). For example, with respect to the first time point, in the framing condition that  
434 did not emphasise 3PP impact on transgressors, children were simply asked: “*So, you might punish*  
435 *some players. How do you think it will feel to do that?*”. Instead, in the framing condition that  
436 emphasised 3PP impact, children were asked: “*So, you might ban some players from the game so*  
437 *they can’t play for quite a while. How do you think it will feel to do that?*” (for details of the  
438 framing manipulation for each time point see Supplementary Information – sections 1.4, 1.7, 1.8).

439 At the end of the experiment, each child was asked a **manipulation check question** to  
440 evaluate the believability of the experimental setting. Specifically, children were questioned about  
441 whether they thought they had actually refereed real internet players during the trials (“*Do you think*  
442 *you really watched games with internet players now?*”).

A

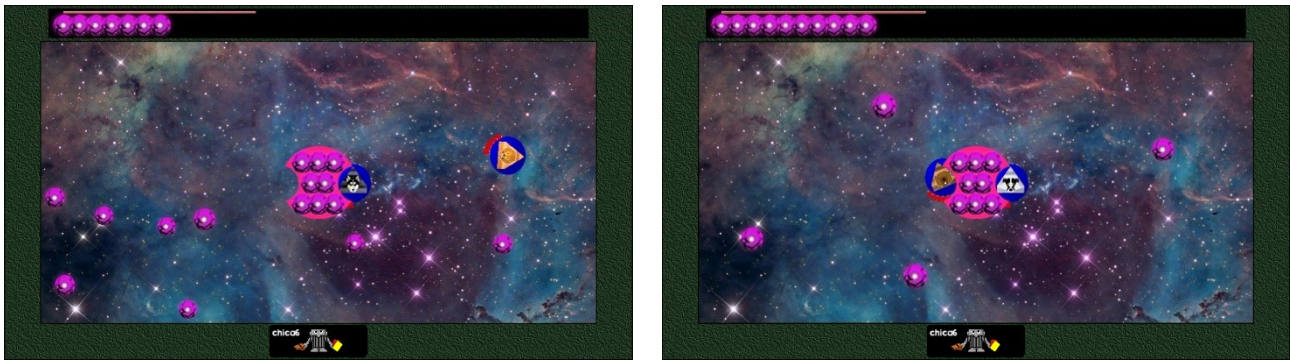


B



C

D



443 **Figure 3. Moral violations shown during the test trials.** (A) *Accidental Unfairness*: player Fox  
 444 takes the majority of the bombs for themselves, leaving only a few for their team-mate, player  
 445 Panda. From the players' voice-overs, the participant can understand that player Fox wanted to  
 446 divide the bombs equally, but failed to do so because they inadvertently used the controls wrongly  
 447 during the bomb distribution. (B) *Failed Intentional Unfairness*: player Kangaroo divides the  
 448 bombs equally between themselves and their team-mate, player Ostrich. From the players' voice-  
 449 overs, the participant can understand that player Kangaroo intended to take the majority of the  
 450 bombs for themselves. However, player Kangaroo did not succeed in their plan, because they used  
 451 the controls wrongly during the bomb distribution. (C) *Accidental Disloyalty*: player Lion fails to  
 452 touch the mega-gem, thus leaving their team-mate, player Dog, stuck in it. From the players' voice-  
 453 overs, the participant can understand that player Lion tried to touch the mega-gem to free player  
 454 Dog from it, but did not succeed in doing so because of their poor spaceship piloting skills. (D)  
 455 *Failed Intentional Disloyalty*: player Beaver touches the mega-gem, thus freeing their team-mate,  
 456 player Badger, who was trapped inside. From the players' voice-overs, the participant can  
 457 understand that player Beaver tried to leave player Badger trapped but, because of their poor  
 458 spaceship piloting skills, ended up involuntarily touching the mega-gem.  
 459

460 ***Analysis Strategy and Statistics***

461 Linear mixed-effects models were used to examine the predictors of punishment severity  
 462 and punishment affective states, with Participants' ID included as a random factor because there  
 463 were multiple data points per individual. All other IVs were included as fixed factors. Notably, we  
 464 included country and gender in our models to explain variance in the DVs, not to test predictions.  
 465 Our aim was not to detect differences between Colombian and Spanish children, but to assess  
 466 whether we would replicate findings previously obtained in Anglo-American or British samples  
 467 (e.g., Cushman et al., 2013; Gummerum & Chu, 2014; Bernhard et al., 2020). Moreover, since there  
 468 was no evidence of gender differences in 3PP in studies using similar paradigms with British  
 469 children (Arini et al., 2021), we did not expect to find gender differences in Colombian and Spanish  
 470 children.

471 Model fits were confirmed by examining diagnostic scatter plots of residuals. All analyses  
472 were conducted in the R programming environment (R version 4.0.2, R Core Team, 2020) with raw  
473 data and code openly available in the OSF repository (<https://osf.io/c9w2a/>). In our models, we  
474 included main effects and, where appropriate to answer our research questions, two-way interaction  
475 effects. However, we did not include any three-way interaction effects, due to concerns of  
476 insufficient power to detect effects because of the small sample size.

477

478

## 479 **Results**

### 480 ***Preliminary Analyses***

481 We firstly evaluated the believability of the experimental setting. Believability was high  
482 (90%, 95% CI [83%, 95%]), indicating that the majority of the children thought they had actually  
483 refereed real internet players during the trials. Since there was almost no variability in this measure,  
484 we excluded believability from our statistical models.

485 Preliminary analyses of the control trial revealed that only 3 out of the 123 participants  
486 identified moral norm violations, when in fact they were shown cases of moral norm conformity  
487 (i.e., both outcomes and intentions were positive for both fairness and loyalty domains). Analyses  
488 presented below therefore exclude the control trial, which served its purpose by demonstrating that  
489 participants could generally distinguish moral norm violations from moral norm conformity.  
490 Therefore, in our statistical models about punishment severity we considered only two within-  
491 subject levels for intentionality: failed intentional transgression and accidental transgression.

### 492 ***Punishment Affective States***

493 As shown in Table 2, linear mixed-effects analyses revealed that there were no significant  
494 temporal changes on children's punishment affective states. Overall, children did not much enjoy  
495 making 3PP decisions: across time points (before, during and after punishment allocation),  
496 children's reported affective states were on average  $M = -.33$ ,  $SD = 2.60$ , which was not

497 significantly different from 0,  $t(122) = -1.42, p = .157$  (Q1 in Table 1; see Figure 4). Additionally,  
 498 whether children were prompted to focus on the impact of 3PP or not had no effect on their  
 499 punishment affective states either (Q2 in Table 1). On the other hand, there was a significant effect  
 500 of country, with Colombian children reporting more negative punishment affective states ( $M = -.92,$   
 501  $SD = 3.55$ ) than Spanish children ( $M = -.03, SD = 2.80$ ).

502

503

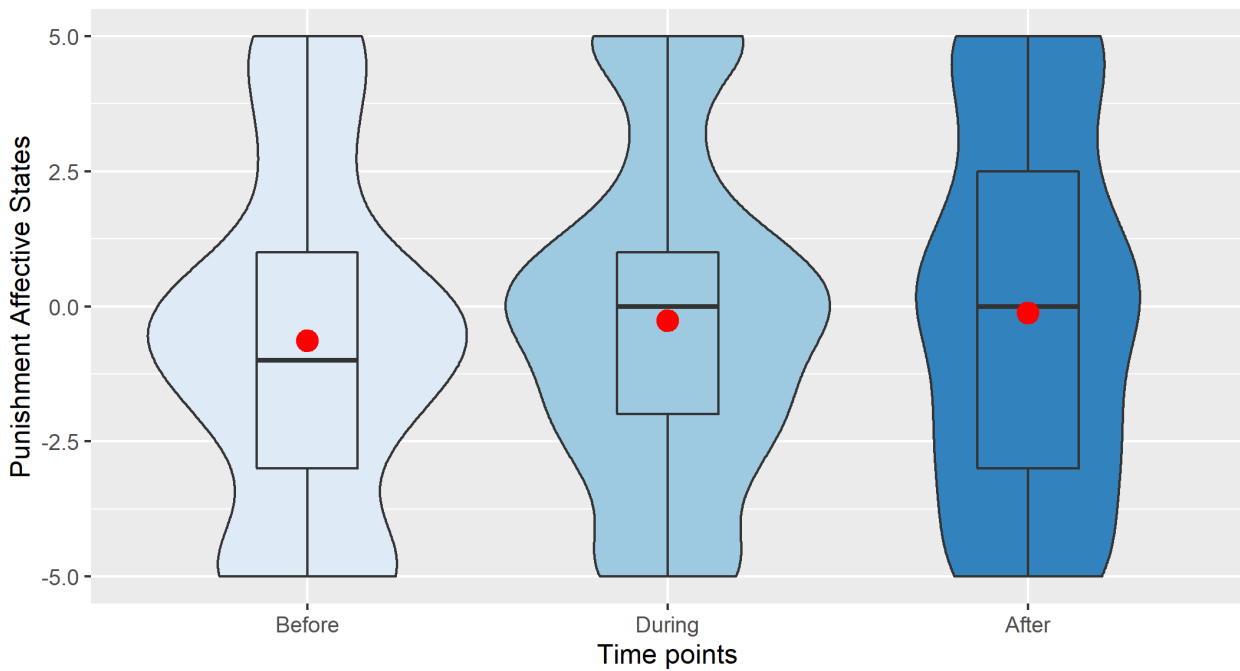
504

505 **Table 2. Modulating factors of punishment affective states.**

Factor	<i>b</i>	$\beta$	95% CI for $\beta$	$\chi^2$	<i>p</i>
Judgement of transgression severity (average)	-.22	-.09	[-.25, .07]	1.25	.263
Age	-.11	-.06	[-.22, .10]	0.53	.467
Gender	.21	.07	[-.24, .37]	0.18	.669
Country	1.17	.38	[.05, .70]	5.07	.024*
Question focus	.19	.06	[-.23, .35]	0.16	.688
Question time				3.93	.140
During vs. Before	.38	.12	[-.05, .29]		
After vs. Before	.51	.16	[.00, .33]		

506 **Note:** \*  $p \leq .050$ . \*\*  $p \leq .010$ . \*\*\*  $p \leq .001$ . For binary variables, the following categories are coded as 1 (and the others as 0):  
 507 gender male, country Spain, question focused on impact, believed the game to be real. An additional categorical variable is  
 508 question time, which is ternary rather than binary (categories are before, during and after, with before used as a reference). The  
 509 continuous predicting factors are age and judgement of transgression severity averaged across trials. Raw model coefficients *b* are  
 510 standardised to produce  $\beta$  and associated 95% confidence interval by normalising by standard deviation of the dependent variable  
 511 in all cases and by the standard deviation of the predicting factor only when it is not categorical (age, judgement of transgression  
 512 severity averaged across trials), meaning categorical  $\beta$  (gender, country, question focus, and question time) is analogous to  
 513 Cohen's *d*. The scale used to measure judgement of transgression severity ranged from -5 to 0, meaning that the more negative the  
 514 values, the harsher/more severe the judgements.

515



516

517 **Figure 4. Punishment affective states across time points (before, during and after punishment**  
 518 **allocation).** Violin plots wrapping boxplots; boxplots showing median and interquartile range,  
 519 outliers, and a large dot for mean value.

520

521

### ***Punishment Severity***

522

523

524

525

526

527

528

529

As shown in Table 3, linear mixed-effects analyses revealed a significant effect of judgement of transgression severity, indicating that the harsher the judgement the more severe the punishment. Additionally, there was a significant effect of age, meaning that the older the children the more lenient their 3PP behaviour towards the transgressors. There was a significant effect of moral domain, with children punishing disloyalty transgressions ( $M = 3.92$ ,  $SD = 1.45$ ) more harshly than unfairness transgressions ( $M = 3.32$ ,  $SD = 1.54$ ). We also found a significant effect of intentionality, meaning that children on average punished failed intentional transgressions ( $M = 3.77$ ,  $SD = 1.45$ ) more harshly than accidental transgressions ( $M = 3.50$ ,  $SD = 1.60$ ).

530

531

532

533

Furthermore, we found a significant interaction between intentionality and age (Q3 in Table 1). Notably, this was not accompanied by an interaction between moral domain and intentionality (Q4 in Table 1). In other words, we found evidence of an outcome-to-intent shift in 3PP behaviour, which occurred in both moral domains in parallel. From a visual interpretation of the data (see

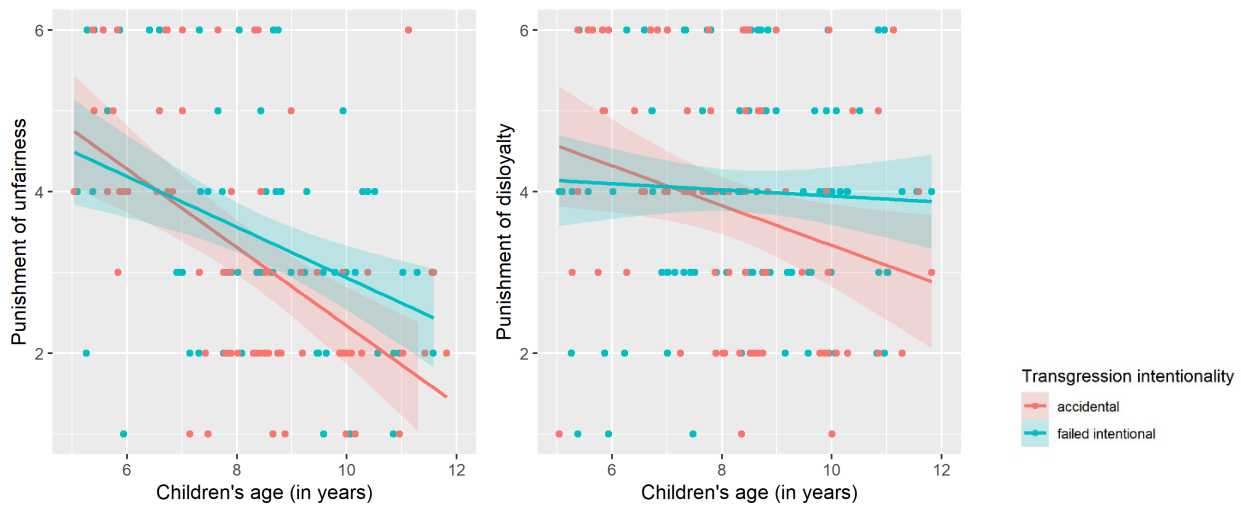
534 Figure 5), it appears that the outcome-to-intent shift occurred around 7 years of age in unfairness  
535 and disloyalty. Namely, children of 7 years of age or younger tended to punish failed intentional  
536 transgressions (disloyalty:  $M = 4.31$ ,  $SD = 1.64$ ; unfairness:  $M = 4.39$ ,  $SD = 1.50$ ) as severely as  
537 accidental transgressions (disloyalty:  $M = 4.61$ ,  $SD = 1.46$ ; unfairness:  $M = 4.78$ ,  $SD = 1.00$ ). In  
538 contrast, children older than 7 tended to punish failed intentional transgressions (disloyalty:  $M =$   
539  $3.92$ ,  $SD = 1.31$ ; unfairness:  $M = 3.15$ ,  $SD = 1.35$ ) more severely than accidental transgressions  
540 (disloyalty:  $M = 3.50$ ,  $SD = 1.48$ ; unfairness:  $M = 2.65$ ,  $SD = 1.43$ ). We additionally discovered a  
541 significant interaction between age and moral domain: punishment severity decreased with  
542 children's increasing age in cases of unfairness, whereas it remained more stable across ages in  
543 cases of disloyalty, see Figure 5.

544 **Table 3. Modulating factors of punishment severity.**

Factor	<i>b</i>	$\beta$	95% CI for $\beta$	$\chi^2$	<i>p</i>
Judgement of transgression severity	-.26	-.26	[-.36, -.15]	22.31	<.001***
Age	-.20	-.22	[-.40, -.05]	25.38	<.001***
Gender	-.01	-.01	[-.26, .24]	0.00	.947
Country	-.12	-.08	[-.35, .19]	0.34	.560
Moral domain	1.31	-.31	[-.56, -.07]	18.27	<.001***
Intentionality	-1.41	.10	[-.12, .32]	9.46	.024*
Moral domain x Intentionality	.12	.08	[-.24, .40]	0.23	.630
Age x Moral domain	-.21	-.24	[-.40, -.09]	9.07	.003**
Age x Intentionality	.19	.21	[.05, .37]	6.84	.009**

545 **Note:** \*  $p \leq .050$ . \*\*  $p \leq .010$ . \*\*\*  $p \leq .001$ . For binary variables, the following categories are coded as 1 (and the others as 0):  
546 gender male, country Spain, believed the game to be real, domain of unfairness, and failed intentional transgression. Raw model  
547 coefficients *b* are standardised to produce  $\beta$  and associated 95% confidence interval by normalising by standard deviation of the  
548 dependent variable in all cases and by the standard deviation of the predicting factor only when it is not categorical (age,  
549 judgement of transgression severity), meaning categorical  $\beta$  (gender, country, moral domain, and intentionality) is analogous to  
550 Cohen's *d*. The scale used to measure judgement of transgression severity ranged from -5 to 0 (the more negative the values, the  
551 harsher the judgements), therefore negative *b* and  $\beta$  coefficients indicate a direct relationship between judgement of transgression  
552 severity and punishment severity (the harsher the judgement, the harsher the punishment).





554

555 **Figure 5. Punishment severity by moral domain (disloyalty, unfairness) and intentionality**  
 556 **(accidental transgression, failed intentional transgression).** Punishment severity is measured on a scale from 1 (no punishment) to 6 (1 day-ban).  
 557  
 558

559

559

560

## Discussion

561

562

563

564

565

566

567

568

569

570

571

572

573

574

By testing 5- to 11-year-old children from Colombia and Spain, the present study has expanded knowledge about cognitive and emotional processes involved in 3PP behaviour – topics that had been investigated so far mainly in samples from Anglo-America or Northwestern Europe (Marshall & McAuliffe, 2022; see also discussions about sampling bias in developmental psychology in Nielsen et al., 2017; Amir & McAuliffe, 2020). We specifically focused on the emotional consequences of implementing 3PP decisions, and on the integration between outcome and intention information in 3PP decision-making, across different moral domains – disloyalty to the group (a group-focused moral domain) and unfairness in resource distribution (an individual-focused moral domain; Graham et al., 2011).

Regarding punishment affective states, we replicated the result that Arini et al. (2021, Study 2) obtained in British children by demonstrating that also Colombian and Spanish children tended not to derive enjoyment from punishing transgressors (although neither did they tend to find it deeply unpleasant). Interestingly, in the present study children’s affective states before 3PP allocation were not more positive than those reported during and after 3PP allocation.



575 Consequently, we can rule out the hypothesis that children had hedonic expectations about 3PP that  
576 did not stand the test of reality. Rather, it is more likely that carrying out 3PP does not usually  
577 evoke much positive emotions in children, in line with Marshall et al. (2021, Supplementary  
578 Information). It is noteworthy that children in our study did not make the same forecasting error  
579 typically committed by adults. Indeed, adults who were asked to predict how they would feel if they  
580 could punish transgressors reported more positive feelings than their counterparts who actually  
581 enacted punishment (Carlsmith et al., 2008). We acknowledge the possibility that, after children had  
582 responded to the first affective question, they simply responded in similar ways on the two  
583 subsequent questions for sake of consistency, rather than because their affective states were really  
584 the same throughout the experiment. However, changes in punishment affective states over time  
585 have been recorded using similar experimental methods, even though children were being  
586 repeatedly asked the same question at different time points (Arini et al., 2023). Therefore, it is  
587 unlikely that our finding (i.e., consistent lack of much enjoyment over time) is a mere artifact due to  
588 the protocol we adopted. Nevertheless, future studies should complement self-reported measures of  
589 punishment affective states with implicit measures of emotional arousal (e.g., skin conductance) to  
590 validate the findings (as in Gummerum et al., 2020). Future research could also test whether  
591 children would be more likely to enjoy 3PP if they were presented with evidence that transgressors  
592 suffered and/or changed moral attitude after punishment, as it has been demonstrated in adults (Eder  
593 et al., 2020; Gollwitzer & Denzler, 2009; Gollwitzer et al., 2011; Funk et al., 2014; Aharoni et al.,  
594 2022).

595 Furthermore, since past research found that children were more likely to report lack of  
596 enjoyment when their 3PP decisions had real rather than pretend consequences (Arini et al., 2021,  
597 Study 2), we decided to investigate whether children's affective states are sensitive to the impact of  
598 3PP on the transgressors. We predicted that inducing children to focus on the impact the 3PP has on  
599 the transgressors while questioning them about their punishment affective states would make them  
600 feel worse, due to feeling responsible for the transgressors' suffering. In fact, it was found that

601 question focus (focus on 3PP impact vs. no focus on 3PP impact) was not a significant predictor of  
602 children's punishment affective states in our experiment. However, since we did not include any  
603 manipulation check to verify whether the wording of the question about punishment affective states  
604 was effective in activating punishment impact representations, this null result is difficult to  
605 interpret. At this stage, we cannot know if our manipulation did not work or if children did not feel  
606 responsible for the transgressor's suffering.

607         We speculate that these affective findings may shed light on children's 3PP motives.  
608 Potential motives guiding punishment are retribution (i.e., desire to make the transgressors suffer in  
609 proportion to the damage they caused as a means to righting past wrongs) and deterrence (i.e.,  
610 desire to make the transgressors learn a moral lesson to prevent them from misbehaving again in the  
611 future; Aharoni et al., 2022). In adults it has been found that manipulating retribution-relevant  
612 information increased participants' punitive tendencies, yet manipulating deterrence-relevant  
613 information did not (Carlsmith et al., 2002; Molho et al., 2022), which suggests that adults are  
614 primarily motivated by retribution. Notably, a piece of information connected to retribution is  
615 punishment severity (Keller et al., 2010, Study 3), which is akin to how we framed 3PP impact in  
616 our experiment. Therefore, the fact that we did not find an effect of 3PP impact on children's  
617 affective states is suggestive that children's 3PP behaviour may not be motivated by retribution.  
618 This would be in accordance with research by Arini et al. (2023), which showed that children not  
619 only endorsed deterrence over retribution (explicit measure of punishment motivation), but also  
620 recalled deterrence messages at higher rates than retribution messages (as an implicit measure of  
621 punishment motivation). Furthermore, there is evidence that children punished transgressors at  
622 higher rates and invested more resources into 3PP when doing so satisfied deterrent goals, in  
623 addition to retributive ones (Marshall et al., 2021; Twardawski & Hilbig, 2020).

624         In our study, we additionally investigated the cognitive integration of outcome and intention  
625 information in actual punishment behaviour in response to ostensibly real moral violations (rather  
626 than in verbal punishment recommendations for hypothetical moral violations, as in most previous

627 studies). We found that children began to punish failed intentional transgressions more severely  
628 than accidental transgressions from around age 7, when the outcome-to-intent shift became  
629 noticeable, with similar patterns across moral domains (unfairness and disloyalty). In contrast,  
630 previous behavioural studies investigating the outcome-to-intent shift in 3PP either found no  
631 evidence of sensitivity to intentions in 4- to 7-year-old children (Bernhard et al., 2020) or found  
632 evidence of this sensitivity only in adulthood (Gummerum & Chu, 2014; Hechler & Kessler, 2022).  
633 Therefore, ours appears to be the first behavioural study to provide evidence of the capability to  
634 integrate outcomes and intentions into 3PP decisions already in childhood. To note, the different  
635 degree of processing demands and stimuli salience between our study and Gummerum & Chu's  
636 (2014) is likely responsible for the striking age difference in the onset of the outcome-to-intent shift  
637 (7 years of age vs. late adolescence). Indeed, the methodology employed by Gummerum & Chu  
638 (2014) may have taxed children's cognitive resources, impeding their capability to integrate  
639 outcomes and intentions into their 3PP decisions (Hilton & Kuhlmeier, 2019; Margoni & Surian,  
640 2020). If confirmed, this would represent further evidence in support of the hypothesis that the  
641 outcome-to-intent shift is affected by the development of cognitive skills (Killen et al., 2011; Zelazo  
642 et al., 1996; for a review see Margoni & Surian, 2016).

643         Additionally, children in our experiment manifested the outcome-to-intent shift in their 3PP  
644 behaviour within the same age range previously observed in the vignette studies on punishment  
645 recommendations (between 5 and 8 years; Baird & Astington, 2004; Cushman et al., 2013; Killen et  
646 al., 2011; Martin et al., 2022; Nobes et al., 2016). This is noteworthy given that individuals often do  
647 not carry out the behaviour they judge appropriate (Blake, 2018; see also discussion in Kenward &  
648 Östh, 2015). A final relevant consideration is that, if 3PP evolved as a mechanism to enforce group  
649 cooperation (Boyd & Richerson, 1992), it makes sense for children to start taking intentions into  
650 account in their 3PP behaviour during middle childhood. It is indeed during this developmental  
651 period that children increasingly engage in social interactions with their peers, and face their first  
652 coordination and bargaining problems (Grueneisen & Tomasello, 2020, 2022). Thus, becoming

653 watchful about clues indicating someone's intention to disregard cooperative norms is likely to be  
654 adaptive, as it would allow the avoidance of apparently unreliable social partners. Consequently, the  
655 cost of developing this ability only in late adolescence would probably be too high (Margoni &  
656 Surian, 2016).

657         Regarding the effect of moral domain on intention sensitivity, children in our experiment  
658 assigned equal weight to intentions in their 3PP decisions across moral domains (unfairness vs.  
659 disloyalty). In contrast, previous studies have shown that adults assign different weights to  
660 intentions depending on the moral domain (with intentions having greater influence in harm than  
661 purity transgressions) (Barrett et al., 2016; Chakroff et al., 2016; Sweetman & Newman, 2020a,  
662 2020b; Young & Saxe, 2011; Young & Tsoi, 2013). Therefore, our findings do not support the  
663 hypothesis that people tend to attribute more importance to intentions in individual-focused  
664 domains (i.e., harm and fairness) than in group-focused domains (i.e., loyalty, authority and purity;  
665 Graham et al., 2011). These contrasting results could be due to differences in the details of the  
666 moral scenarios or to developmental differences. Differences in intention sensitivity may only be  
667 detectable when comparing harm and purity transgressions, and may not extend to other moral  
668 domains, or all exemplars within them. Alternatively, children might not have fully developed the  
669 ability to differentiate the weight of intentions across different moral domains. If confirmed, this  
670 would suggest that moral decision-making becomes more domain-specific with age. Future studies  
671 should therefore discern between alternative explanations by investigating the developmental  
672 trajectory of intention sensitivity from childhood to adulthood across a broader range of moral  
673 domains and scenarios within them.

674         However, it is worth noting our exploratory analyses about how punishment of different  
675 moral transgressions changes across development. In this experiment, Colombian and Spanish  
676 children punished disloyalty more harshly than unfairness. Moreover, their 3PP severity of  
677 disloyalty tended to remain stable across ages, while 3PP severity of unfairness decreased as  
678 children got older. It could be argued that 3PP severity of disloyalty remained stable throughout

679 development because the disloyalty scenario was less cognitively demanding than the unfairness  
680 scenario. However, this would not explain why British children in the same age range reacted to the  
681 view of the same moral scenarios in quite different ways (Arini et al., 2021). When British children  
682 were tested on a paradigm that closely resembled the one Colombian and Spanish children were  
683 confronted with, their 3PP severity was comparable across moral domains and decreased with an  
684 age-dependent pattern for disloyalty and unfairness alike (Arini et al., 2021, Study 2). When instead  
685 British children were tested on a paradigm that allowed them to use 3PP not only to make the  
686 transgressor pay for their action but also to equalise the resource unbalance between victim and  
687 transgressor (an arguably cognitively demanding task), 3PP severity of unfairness remained steadily  
688 high across ages, while 3PP severity of disloyalty decreased as children got older (Arini et al., 2021,  
689 Study 1). Considering this evidence, even though we cannot rule out that differences in cognitive  
690 demands between moral scenarios played a role in children's 3PP decisions, we deem them unlikely  
691 to explain our pattern of results.

692 To sum up, if we consider punishment severity as a proxy of the importance attributed to a  
693 specific moral domain, Colombian and Spanish children were more concerned about disloyalty than  
694 unfairness, whereas British children were either equally concerned about the two, or more  
695 concerned about unfairness than disloyalty (Arini et al., 2021). Given that the culture in Spain and  
696 even more in Colombia is more collectivist than in the UK (Hofstede, 2001; see also Kryś et al.,  
697 2022; Uskul et al., 2023), these findings are in line with research conducted in adults suggesting  
698 that collectivism may be associated with higher concerns about group- than individual-focused  
699 moral domains (Graham et al., 2011; Triandis, 1989). Crucially, when differences between moral  
700 concerns were detected within a sample, either in the present experiment or in Arini et al. (2021),  
701 they tended to increase with development. In other words, the longer children were exposed to the  
702 specific moral system of their own socio-cultural environment, the more their moral concerns  
703 became selective towards the moral domain deemed central in said environment, thus mirroring  
704 adults' moral concerns – a pattern consistent with cultural learning processes. This is an important

705 testing ground for moral foundations theory's claim that moral development is driven by cultural  
706 learning (Graham et al., 2013). Although the results are suggestive, this interpretation warrants  
707 caution. The samples of children in both the present study and Arini et al. (2021) were not  
708 necessarily representative of the respective national populations. It follows that their punishment  
709 behaviour may reflect local norms in their specific environment (e.g., school, neighbourhood) rather  
710 than collectivist or individualistic tendencies in their countries (Colombia, Spain, UK). However, if  
711 future research confirmed this preliminary evidence, it would provide insight into the complex  
712 relationship between culturally-salient moral norms and the development and variation of children's  
713 3PP behaviour across societies. Such studies would also benefit from taking a gender perspective  
714 since there is evidence of gender differences in individualism and collectivism (Dabiriyani Tehrani  
715 & Yamini, 2022), which may affect moral concerns towards specific norm violations.

716 Finally, it is important to acknowledge that a limitation of our study due to logistic  
717 constraints is the relatively small size, which might have prevented the detection of effects when  
718 they were in fact present because of lack of statistical power. This shortcoming might have also  
719 created issues of reliability for the effects that were indeed detected. Therefore, the current evidence  
720 should be regarded as preliminary, and future studies should aim at replicating our results in a larger  
721 sample. Moreover, there were differences in size, gender distribution, mean age and  
722 representativeness between the Colombian and the Spanish samples. The Colombian sample was  
723 smaller, with a higher proportion of male children, and its mean age was a whole year younger  
724 compared to the Spanish sample. Additionally, all Colombian participants came from the same city  
725 and attended the same school, whereas Spanish participants were recruited from four different  
726 schools located in two different cities. However, our choice to recruit children from Colombia and  
727 Spain was not motivated by the desire to detect cultural differences between these two samples.  
728 Rather, we wanted to broaden representation in developmental psychology (Nielsen et al., 2017;  
729 Amir & McAuliffe, 2020) and test the generalisability of findings about children's 3PP previously  
730 obtained in Anglo-American or Northwestern European samples (reviewed in Marshall &

731 McAuliffe, 2022). Another limitation of the current study is that it was conducted in one specific  
732 experimental setting – a computer-mediated paradigm – whose generalisability to real-world  
733 situations has only recently been tested (Arini et al., 2023). However, computer games represent a  
734 real social world that children already inhabit, experience and react to norm violations within (Kou  
735 et al., 2017), thus the ecological validity of the present experimental paradigm is expected to be  
736 high. A further weakness of our study is that we employed only one behavioural exemplar for each  
737 moral domain, therefore future studies would benefit from using more than one example of  
738 behaviour per each type of moral domain. Finally, we acknowledge the lack of order balancing for  
739 the control trial (only test trials were counterbalanced). This design choice was motivated by the  
740 need to initiate the refereeing sessions with the same baseline condition for all the participants  
741 (similarly to what has been done in e.g. Twardawski & Hilbig, 2020 and Arini et al., 2021), but it  
742 would be beneficial if future studies adopted a fully counterbalanced design as a robustness check.

743         In conclusion, the present study has deepened the understanding of cognitive and emotional  
744 processes playing a crucial role in children’s moral development. To our knowledge, this has been  
745 the first study to provide evidence of the outcome-to-intent shift in 3PP behaviour during middle  
746 childhood. More specifically, children began to attribute higher importance to intentions over  
747 outcomes in 3PP behaviour, across different moral domains, around 7 years of age, in line with  
748 findings about the outcome-to-intent shift in punishment recommendations (Baird & Astington,  
749 2004; Cushman et al., 2013; Killen et al., 2011; Martin et al., 2022; Nobes et al., 2016). We also  
750 found that children in our study did not derive much enjoyment from enacting 3PP, in accordance  
751 with previous literature (Arini et al., 2021, Study 2; Marshall et al., 2021, Supplementary  
752 Information), nor did they anticipate to feel much enjoyment. We also discovered interesting cross-  
753 cultural differences: Colombian and Spanish children punished disloyalty more severely than  
754 unfairness, in contrast with the behavioural patterns observed in British children, whose 3PP  
755 severity of unfairness was either higher or equal to that of disloyalty (Arini et al., 2021). Since  
756 different cultures privilege different moral domains (Graham et al., 2011; Triandis, 1989), further

757 studies are needed at the intersection between developmental psychology and cognitive  
758 anthropology in order to shed light on moral development from a cross-cultural perspective. This  
759 would enable a more fine-grained distinction between universal and culture-specific developmental  
760 patterns of punishment behaviour and affective states, ultimately enriching understanding about  
761 proximate and evolutionary causes of our socio-moral behaviour.

## 762 **Acknowledgements**

763 We thank Morag MacLean and Michaela Gummerum for their precious feedback on an  
764 early draft of the manuscript, all the school staff who supported this project (particularly Gina  
765 Piñeros Romero), and the families who took part in it. Additionally, we thank Daniela Rosas  
766 Moreno for assistance in the development of the methodology, Juan Sebastian Nassar Pereira and  
767 Juan Jaccobo Garzon Martelo for assistance in data collection, and Vanessa Anahi Merlos-Fuentes  
768 for subtitling the video rendition of the refereeing sessions.

## 769 **Funding Sources**

770 This work was supported by the Nigel Groome Studentship and Santander Research  
771 Scholarship (both awarded to Rhea L. Arini) and internal research funding from the Faculty of  
772 Social Sciences, Universidad de los Andes, Colombia (awarded to Gordon P. D. Ingram).

## 773 **Declarations of Interest**

774 None.

## 775 **References**

776 Aharoni, E., Simpson, D., Nahmias, E., & Gollwitzer, M. (2022). A painful message: Testing the  
777 effects of suffering and understanding on punishment judgments. *Zeitschrift für*  
778 *Psychologie*. <https://doi.org/10.1027/2151-2604/a000460>



- 779 Amir, D., & McAuliffe, K. (2020). Cross-cultural, developmental psychology: Integrating  
780 approaches and key insights. *Evolution and Human Behavior*, 41(5), 430-444.  
781 <https://doi.org/10.1016/j.evolhumbehav.2020.06.006>
- 782 Arini, R. L., Mahmood, M., Aljure, J. B., Ingram, G. P., Wiggs, L., & Kenward, B. (2023).  
783 Children endorse deterrence motivations for third-party punishment but derive higher  
784 enjoyment from compensating victims. *Journal of Experimental Child Psychology*, 230,  
785 105630. <https://doi.org/10.1016/j.jecp.2023.105630>
- 786 Arini, R. L., Wiggs, L., & Kenward, B. (2021). Moral duty and equalization concerns motivate  
787 children's third-party punishment. *Developmental Psychology*, 57(8), 1325-1341.  
788 <https://doi.org/10.1037/dev0001191>
- 789 Baird, J. A., & Astington, J. W. (2004). The role of mental state understanding in the development  
790 of moral cognition and moral action. *New Directions for Child and Adolescent*  
791 *Development*, 2004(103), 37-49. <https://doi.org/10.1002/cd.96>
- 792 Barrett, H. C., Bolyanatz, A., Crittenden, A. N., Fessler, D. M., Fitzpatrick, S., Gurven, M., . . .  
793 Pisor, A. (2016). Small-scale societies exhibit fundamental variation in the role of intentions  
794 in moral judgment. *Proceedings of the National Academy of Sciences*, 113(17), 4688-4693.  
795 <https://doi.org/10.1073/pnas.1522070113>
- 796 Bernhard, R. M., Martin, J. W., & Warneken, F. (2020). Why do children punish? Fair outcomes  
797 matter more than intent in children's second- and third-party punishment. *Journal of*  
798 *Experimental Child Psychology*, 200, 104909. <https://doi.org/10.1016/j.jecp.2020.104909>
- 799 Blake, P. R. (2018). Giving what one should: Explanations for the knowledge-behavior gap for  
800 altruistic giving. *Current Opinion in Psychology*, 20, 1-5.  
801 <https://doi.org/10.1016/j.copsyc.2017.07.041>
- 802 Boyd, R., & Richerson, P. J. (1992). Punishment allows the evolution of cooperation (or anything  
803 else) in sizable groups. *Ethology and Sociobiology*, 13(3), 171-195.  
804 [https://doi.org/10.1016/0162-3095\(92\)90032-Y](https://doi.org/10.1016/0162-3095(92)90032-Y)

805 Bueno-Guerra, N., Leiva, D., Colell, M., & Call, J. (2016). Do sex and age affect strategic behavior  
806 and inequity aversion in children?. *Journal of Experimental Child Psychology*, *150*, 285-  
807 300. <https://doi.org/10.1016/j.jecp.2016.05.011>

808 Carlsmith, K. M., Darley, J. M., & Robinson, P. H. (2002). Why do we punish? Deterrence and just  
809 deserts as motives for punishment. *Journal of Personality and Social Psychology*, *83*(2),  
810 284-299. <https://doi.org/10.1037/0022-3514.83.2.284>

811 Carlsmith, K. M., Wilson, T. D., & Gilbert, D. T. (2008). The paradoxical consequences of revenge.  
812 *Journal of Personality and Social Psychology*, *95*(6), 1316.  
813 <https://doi.org/10.1037/a0012165>

814 Chakroff, A., Dungan, J., Koster-Hale, J., Brown, A., Saxe, R., & Young, L. (2016). When minds  
815 matter for moral judgment: intent information is neurally encoded for harmful but not  
816 impure acts. *Social Cognitive and Affective Neuroscience*, *11*(3), 476-484.  
817 <https://doi.org/10.1093/scan/nsv131>

818 Chernyak, N., & Sobel, D. M. (2016). “But he didn’t mean to do it”: Preschoolers correct  
819 punishments imposed on accidental transgressors. *Cognitive Development*, *39*, 13-20.  
820 <https://doi.org/10.1016/j.cogdev.2016.03.002>

821 Curry, O. S., Chesters, M. J., & Van Lissa, C. J. (2019). Mapping morality with a compass: Testing  
822 the theory of ‘morality-as-cooperation’ with a new questionnaire. *Journal of Research in*  
823 *Personality*, *78*, 106-124. <https://doi.org/10.1016/j.jrp.2018.10.008>

824 Cushman, F. (2008). Crime and punishment: Distinguishing the roles of causal and intentional  
825 analyses in moral judgment. *Cognition*, *108*(2), 353-380.  
826 <https://doi.org/10.1016/j.cognition.2008.03.006>

827 Cushman, F., Sheketoff, R., Wharton, S., & Carey, S. (2013). The development of intent-based  
828 moral judgment. *Cognition*, *127*(1), 6-21. <https://doi.org/10.1016/j.cognition.2012.11.008>

- 829 Dabiriyani Tehrani, H., & Yamini, S. (2022). Gender differences concerning the horizontal and  
830 vertical individualism and collectivism: A meta-analysis. *Psychological Studies*, 67(1), 11-  
831 27. <https://doi.org/10.1007/s12646-022-00638-x>
- 832 De Quervain, D. J. F., Fischbacher, U., Treyer, V., Schellhammer, M., Schnyder, U., Buck, A., &  
833 Fehr, E. (2004). The neural basis of altruistic punishment. *Science*, 305(5688), 1254-1258.  
834 <https://doi.org/10.1126/science.1100735>
- 835 Eder, A. B., Mitschke, V., & Gollwitzer, M. (2020). What stops revenge taking? Effects of  
836 observed emotional reactions on revenge seeking. *Aggressive Behavior*, 46(4), 305-316.  
837 <https://doi.org/10.1002/ab.21890>
- 838 Fehr, E., & Gächter, S. (2002). Altruistic punishment in humans. *Nature*, 415(6868), 137-140.  
839 <https://doi.org/10.1038/415137a>
- 840 Funk, F., McGeer, V., & Gollwitzer, M. (2014). Get the message: Punishment is satisfying if the  
841 transgressor responds to its communicative intent. *Personality and Social Psychology*  
842 *Bulletin*, 40(8), 986-997. <https://doi.org/10.1177/0146167214533130>
- 843 Ginther, M. R., Hartsough, L. E., & Marois, R. (2022). Moral outrage drives the interaction of harm  
844 and culpable intent in third-party punishment decisions. *Emotion*, 22(4), 795.  
845 <https://doi.org/10.1037/emo0000950>
- 846 Gollwitzer, M., & Denzler, M. (2009). What makes revenge sweet: Seeing the offender suffer or  
847 delivering a message? *Journal of Experimental Social Psychology*, 45(4), 840-844.  
848 <https://doi.org/10.1016/j.jesp.2009.03.001>
- 849 Gollwitzer, M., Meder, M., & Schmitt, M. (2011). What gives victims satisfaction when they seek  
850 revenge? *European Journal of Social Psychology*, 41(3), 364-374.  
851 <https://doi.org/10.1002/ejsp.782>
- 852 Gonzalez-Gadea, M. L., Dominguez, A., & Petroni, A. (2022). Decisions and mechanisms of  
853 intergroup bias in children's third-party punishment. *Social Development*, 31(4), 1194-1210.  
854 <https://doi.org/10.1111/sode.12608>.

- 855 Graham, J., Haidt, J., Koleva, S., Motyl, M., Iyer, R., Wojcik, S. P., & Ditto, P. H. (2013). Moral  
856 foundations theory: The pragmatic validity of moral pluralism. In *Advances in Experimental*  
857 *Social Psychology* (Vol. 47, pp. 55-130). Academic Press.
- 858 Graham, J., Nosek, B. A., Haidt, J., Iyer, R., Koleva, S., & Ditto, P. H. (2011). Mapping the moral  
859 domain. *Journal of Personality and Social Psychology*, 101(2), 366.  
860 <https://doi.org/10.1037/a0021847>
- 861 Gummerum, M., & Chu, M. T. (2014). Outcomes and intentions in children's, adolescents', and  
862 adults' second-and third-party punishment behavior. *Cognition*, 133(1), 97-103.  
863 <https://doi.org/10.1016/j.cognition.2014.06.001>
- 864 Grueneisen, S., & Tomasello, M. (2020). The development of coordination via joint expectations  
865 for shared benefits. *Developmental Psychology*, 56(6), 1149-1156.  
866 <https://doi.org/10.1037/dev0000936>
- 867 Grueneisen, S., & Tomasello, M. (2022). How fairness and dominance guide young children's  
868 bargaining decisions. *Child Development*, 93(5), 1318-1333.  
869 <https://doi.org/10.1111/cdev.13757>
- 870 Gummerum, M., López-Pérez, B., Van Dijk, E., & Van Dillen, L. F. (2022). Ire and punishment:  
871 incidental anger and costly punishment in children, adolescents, and adults. *Journal of*  
872 *Experimental Child Psychology*, 218, 105376. <https://doi.org/10.1016/j.jecp.2022.105376>
- 873 Gummerum, M., López-Pérez, B., Van Dijk, E., & Van Dillen, L. F. (2020). When punishment is  
874 emotion-driven: Children's, adolescents', and adults' costly punishment of unfair  
875 allocations. *Social Development*, 29(1), 126-142. <https://doi.org/10.1111/sode.12387>
- 876 Gummerum, M., Takezawa, M., & Keller, M. (2009). The influence of social category and  
877 reciprocity on adults' and children's altruistic behavior. *Evolutionary Psychology*, 7(2),  
878 147470490900700. <https://doi.org/10.1177/147470490900700212>

879 Grođlu, B., van den Bos, W., & Crone, E. A. (2009). Fairness considerations: increasing  
880 understanding of intentionality during adolescence. *Journal of Experimental Child*  
881 *Psychology*, 104(4), 398-409. <https://doi.org/10.1016/j.jecp.2009.07.002>

882 Grođlu, B., van den Bos, W., van Dijk, E., Rombouts, S. A., & Crone, E. A. (2011). Dissociable  
883 brain networks involved in development of fairness considerations: Understanding  
884 intentionality behind unfairness. *Neuroimage*, 57(2), 634-641.  
885 <https://doi.org/10.1016/j.neuroimage.2011.04.032>

886 Hamlin, J. K. (2013). Failed attempts to help and harm: Intention versus outcome in preverbal  
887 infants' social evaluations. *Cognition*, 128(3), 451-474.  
888 <https://doi.org/10.1016/j.cognition.2013.04.004>

889 Hamlin, J. K., Wynn, K., Bloom, P., & Mahajan, N. (2011). How infants and toddlers react to  
890 antisocial others. *Proceedings of the National Academy of Sciences of the United States of*  
891 *America*, 108(50), 19931-19936. <https://doi.org/10.1073/pnas.1110306108>

892 Hartsough, L. E., Ginther, M. R., & Marois, R. (2020). Distinct affective responses to second-and  
893 third-party norm violations. *Acta Psychologica*, 205, 103060.  
894 <https://doi.org/10.1016/j.actpsy.2020.103060>

895 Hechler, S., & Kessler, T. (2022). The importance of unfair intentions and outcome inequality for  
896 punishment by third parties and victims. *Zeitschrift fr Psychologie*.  
897 <https://doi.org/10.1027/2151-2604/a000458>

898 Helwig, C. C., Zelazo, P. D., & Wilson, M. (2001). Children's judgments of psychological harm in  
899 normal and noncanonical situations. *Child Development*, 72(1), 66-81.  
900 <https://doi.org/10.1111/1467-8624.00266>

901 Hilton, B. C., & Kuhlmeier, V. A. (2019). Intention attribution and the development of moral  
902 evaluation. *Frontiers in Psychology*, 9, 2663. <https://doi.org/10.3389/fpsyg.2018.02663>

903 Hofstede, G. (2001). *Culture's consequences: Comparing values, behaviors, institutions and*  
904 *organizations across nations*. Sage Publications.

- 905 House, B. R., Kanngiesser, P., Barrett, H. C., Yilmaz, S., Smith, A. M., Sebastian-Enesco, C., ... &  
906 Silk, J. B. (2020). Social norms and cultural diversity in the development of third-party  
907 punishment. *Proceedings of the Royal Society B*, 287(1925), 20192794.  
908 <https://doi.org/10.1098/rspb.2019.2794>
- 909 Ingram, G. P., & Moreno-Romero, C. (2021). Dual-process theories, cognitive decoupling and the  
910 outcome-to-intent shift: A developmental perspective on evolutionary ethics. In: De Smedt,  
911 J., De Cruz, H. (eds) *Empirically Engaged Evolutionary Ethics* (pp. 17-40). Synthese  
912 Library, vol 437. Springer, Cham. [https://doi.org/10.1007/978-3-030-68802-8\\_2](https://doi.org/10.1007/978-3-030-68802-8_2)
- 913 Janoff-Bulman, R., Sheikh, S., & Hepp, S. (2009). Proscriptive versus prescriptive morality: two  
914 faces of moral regulation. *Journal of Personality and Social Psychology*, 96(3), 521.  
915 <http://dx.doi.org/10.1037/a0013779>
- 916 Jaroslawska, A. J., McCormack, T., Burns, P., & Caruso, E. M. (2020). Outcomes versus intentions  
917 in fairness-related decision making: School-aged children's decisions are just like those of  
918 adults. *Journal of Experimental Child Psychology*, 189, 104704.  
919 <https://doi.org/10.1016/j.jecp.2019.104704>
- 920 Jensen, K. (2010). Punishment and spite, the dark side of cooperation. *Philosophical Transactions*  
921 *of the Royal Society B: Biological Sciences*, 365(1553), 2635-2650.  
922 <https://doi.org/10.1098/rstb.2010.0146>
- 923 Jordan, J. J., McAuliffe, K., & Warneken, F. (2014). Development of in-group favoritism in  
924 children's third-party punishment of selfishness. *Proceedings of the National Academy of*  
925 *Sciences of the United States of America*, 111(35), 12710-12715.  
926 <https://doi.org/10.1073/pnas.1402280111>
- 927 Keller, L. B., Oswald, M. E., Stucki, I., & Gollwitzer, M. (2010). A closer look at an eye for an eye:  
928 Laypersons' punishment decisions are primarily driven by retributive motives. *Social*  
929 *Justice Research*, 23(2-3), 99-116. <https://doi.org/10.1007/s11211-010-0113-4>

- 930 Kenward, B., & Östh, T. (2015). Five-year-olds punish antisocial adults. *Aggressive Behavior*,  
931 41(5). <https://doi.org/10.1002/AB.21568>
- 932 Killen, M., Mulvey, K. L., Richardson, C., Jampol, N., & Woodward, A. (2011). The accidental  
933 transgressor: Testing theory of mind and morality knowledge in young  
934 children. *Cognition*, 119, 197-215. <https://doi.org/10.1016/j.cognition.2011.01.006>
- 935 Kou, Y., Johansson, M., & Verhagen, H. (2017). In *Prosocial behavior in an online game*  
936 *community: An ethnographic study* (pp. 1-6). Association for Computing Machinery.  
937 <https://doi.org/10.1145/3102071.3102078>
- 938 Krys, K., Vignoles, V. L., de Almeida, I., & Uchida, Y. (2022). Outside the “Cultural Binary”:  
939 Understanding Why Latin American Collectivist Societies Foster Independent  
940 Selves. *Perspectives on Psychological Science*, 17(4), 1166-  
941 1187. <https://doi.org/10.1177/17456916211029632>
- 942 Lee, Y. E., & Warneken, F. (2022). Does third-party punishment in children aim at equality?.  
943 *Developmental Psychology*, 58(5), 866. <https://doi.org/10.1037/dev0001331>
- 944 Lotz, S., Okimoto, T. G., Schlösser, T., & Fetchenhauer, D. (2011). Punitive versus compensatory  
945 reactions to injustice: Emotional antecedents to third-party interventions. *Journal of*  
946 *Experimental Social Psychology*, 47(2), 477-480. <https://doi.org/10.1016/j.jesp.2010.10.004>
- 947 Margoni, F., & Surian, L. (2016). Explaining the U-shaped development of intent-based moral  
948 judgments. *Frontiers in Psychology*, 7, 171613. <https://doi.org/10.3389/fpsyg.2016.00219>
- 949 Margoni, F., & Surian, L. (2017). Children’s intention-based moral judgments of helping  
950 agents. *Cognitive Development*, 41, 46-64. <https://doi.org/10.1016/j.cogdev.2016.12.001>
- 951 Margoni, F., & Surian, L. (2020). Conceptual continuity in the development of intent-based moral  
952 judgment. *Journal of Experimental Child Psychology*, 194, 104812.  
953 <https://doi.org/10.1016/j.jecp.2020.104812>



954 Marlowe, F. W., Berbesque, J. C., Barr, A., Barrett, C., Bolyanatz, A., Cardenas, J. C., ... & Tracer,  
955 D. (2008). More ‘altruistic’ punishment in larger societies. *Proceedings of the Royal Society*  
956 *B: Biological Sciences*, 275(1634), 587-592. <https://doi.org/10.1098/rspb.2007.1517>

957 Marshall, J., & McAuliffe, K. (2022). Children as assessors and agents of third-party  
958 punishment. *Nature Reviews Psychology*, 1(6), 334-344. [https://doi.org/10.1038/s44159-](https://doi.org/10.1038/s44159-022-00046-y)  
959 [022-00046-y](https://doi.org/10.1038/s44159-022-00046-y)

960 Marshall, J., Yudkin, D. A., & Crockett, M. J. (2021). Children punish third parties to satisfy both  
961 consequentialist and retributive motives. *Nature Human Behaviour*, 5(3), 361-368.  
962 <https://doi.org/10.1038/s41562-020-00975-9>

963 Martin, J. W., Leddy, K., Young, L., & McAuliffe, K. (2022). An earlier role for intent in children’s  
964 partner choice versus punishment. *Journal of Experimental Psychology: General*, 151(3),  
965 597. <https://doi.org/10.1037/xge0001093>

966 McAuliffe, K., Jordan, J. J., & Warneken, F. (2015). Costly third-party punishment in young  
967 children. *Cognition*, 134, 1-10. <https://doi.org/10.1016/j.cognition.2014.08.013>

968 Mendes, N., Steinbeis, N., Bueno-Guerra, N., Call, J., & Singer, T. (2018). Preschool children and  
969 chimpanzees incur costs to watch punishment of antisocial others. *Nature Human Behaviour*  
970 2(1), 45–51. <https://doi.org/10.1038/s41562-017-0264-5>

971 Molho, C., Twardawski, M., & Fan, L. (2022). What motivates direct and indirect punishment?  
972 Extending the “intuitive retributivism” hypothesis. *Zeitschrift für Psychologie*, 230(2), 84-  
973 93. <https://doi.org/10.1027/2151-2604/a000455>

974 Molho, C., Tybur, J. M., Güler, E., Balliet, D., & Hofmann, W. (2017). Disgust and anger relate to  
975 different aggressive responses to moral violations. *Psychological Science*, 28(5), 609-619.  
976 <https://doi.org/10.1177/0956797617692000>

977 Nelson, S. A. (1980). Factors influencing young children’s use of motives and outcomes as moral  
978 criteria. *Child Development*, 823-829. <https://doi.org/10.2307/1129470>



- 979 Nielsen, M., Haun, D., Kärtner, J., & Legare, C. H. (2017). The persistent sampling bias in  
980 developmental psychology: A call to action. *Journal of Experimental Child Psychology*,  
981 *162*, 31-38. <https://doi.org/10.1016/j.jecp.2017.04.017>
- 982 Nobes, G., Panagiotaki, G., & Bartholomew, K. J. (2016). The influence of intention, outcome and  
983 question-wording on children's and adults' moral judgments. *Cognition*, *157*, 190-204.  
984 <https://doi.org/10.1016/j.cognition.2016.08.019>
- 985 Nobes, G., Panagiotaki, G., & Pawson, C. (2009). The influence of negligence, intention, and  
986 outcome on children's moral judgments. *Journal of Experimental Child Psychology*, *104*(4),  
987 382-397. <https://doi.org/10.1016/j.jecp.2009.08.001>
- 988 Pelligra, V., Isoni, A., Fadda, R., & Doneddu, G. (2015). Theory of mind, perceived intentions and  
989 reciprocal behaviour: Evidence from individuals with Autism Spectrum Disorder. *Journal of*  
990 *Economic Psychology*, *49*, 95-107. <https://doi.org/10.1016/j.joep.2015.05.001>
- 991 Pfattheicher, S., Sassenrath, C., & Keller, J. (2019). Compassion magnifies third-party  
992 punishment. *Journal of Personality and Social Psychology*, *117*(1), 124.  
993 <https://doi.org/10.1037/pspi0000165>
- 994 R Core Team (2020). *R: A language and environment for statistical computing*. R Foundation for  
995 *Statistical Computing*. Vienna, Austria. <https://www.R-project.org/>
- 996 Raihani, N. J., & Bshary, R. (2019). Punishment: one tool, many uses. *Evolutionary Human*  
997 *Sciences*, *1*, e12. <https://doi.org/10.1017/ehs.2019.12>
- 998 Riedl, K., Jensen, K., Call, J., & Tomasello, M. (2015). Restorative justice in children. *Current*  
999 *Biology*, *25*(13), 1731-1735. <https://doi.org/10.1016/j.cub.2015.05.014>
- 1000 Salali, G. D., Juda, M., & Henrich, J. (2015). Transmission and development of costly punishment  
1001 in children. *Evolution and Human Behavior*, *36*(2), 86-94.  
1002 <https://doi.org/10.1016/j.evolhumbehav.2014.09.004>

1003 Strobel, A., Zimmermann, J., Schmitz, A., Reuter, M., Lis, S., Windmann, S., & Kirsch, P. (2011).  
1004 Beyond revenge: neural and genetic bases of altruistic punishment. *Neuroimage*, 54(1), 671-  
1005 680. <https://doi.org/10.1016/j.neuroimage.2010.07.051>

1006 Sutter, M. (2007). Outcomes versus intentions: On the nature of fair behavior and its development  
1007 with age. *Journal of Economic Psychology*, 28(1), 69-78.  
1008 <https://doi.org/10.1016/j.joep.2006.09.001>

1009 Sweetman, J., & Newman, G. A. (2020a). Replicating different roles of intent across moral  
1010 domains. *Royal Society Open Science*, 7(5), 190808. <https://doi.org/10.1098/rsos.190808>

1011 Sweetman, J., & Newman, G. A. (2020b). Attentional efficiency does not explain the mental state×  
1012 domain effect. *Plos One*, 15(6), e0234500. <https://doi.org/10.1371/journal.pone.0234500>

1013 Tomasello, M., Carpenter, M., Call, J., Behne, T., & Moll, H. (2005). Understanding and sharing  
1014 intentions: The origins of cultural cognition. *Behavioral and Brain Sciences*, 28(5), 675-691.  
1015 <https://doi.org/10.1017/S0140525X05000129>

1016 Triandis, H. C. (1989). The self and social behavior in differing cultural contexts. *Psychological*  
1017 *Review*, 96(3), 506. <https://doi.org/10.1037/0033-295X.96.3.506>

1018 Twardawski, M., & Hilbig, B. E. (2020). The motivational basis of third-party punishment in  
1019 children. *PLoS One*, 15(11), e0241919. <https://doi.org/10.1371/journal.pone.0241919>

1020 Tybur, J. M., Molho, C., Cakmak, B., Cruz, T. D., Singh, G. D., & Zwicker, M. (2020). Disgust,  
1021 anger, and aggression: Further tests of the equivalence of moral emotions. *Collabra:*  
1022 *Psychology*, 6(1), 34. <https://doi.org/10.1525/collabra.349>

1023 Uskul, A. K., Kirchner-Häusler, A., Vignoles, V. L., Rodriguez-Bailón, R., Castillo, V. A., Cross,  
1024 S. E., Yalçın, M. G., Harb, C., Husnu, S., Ishii, K., Jin, S., Karamaouna, P., Kafetsios, K.,  
1025 Kateri, E., Matamoros-Lima, J., Liu, D., Miniesy, R., Na, J., Özkan, Z., . . . Uchida, Y.  
1026 (2023). Neither Eastern nor Western: Patterns of independence and interdependence in  
1027 Mediterranean societies. *Journal of Personality and Social Psychology*, 125(3), 471-  
1028 495. <https://doi.org/10.1037/pspa0000342>

- 1029 Van de Vondervoort, J. W., & Hamlin, J. K. (2018). Preschoolers focus on others' intentions when  
1030 forming sociomoral judgments. *Frontiers in Psychology*, 9.  
1031 <https://doi.org/10.3389/fpsyg.2018.01851>
- 1032 van den Bos, W., van Dijk, E., & Crone, E. A. (2012). Learning whom to trust in repeated social  
1033 interactions: A developmental perspective. *Group Processes & Intergroup Relations*, 15(2),  
1034 243-256. <https://doi.org/10.1177/1368430211418698>
- 1035 Wittig, M., Jensen, K., & Tomasello, M. (2013). Five-year-olds understand fair as equal in a mini-  
1036 ultimatum game. *Journal of Experimental Child Psychology*, 116(2), 324-337.  
1037 <https://doi.org/10.1016/j.jecp.2013.06.004>
- 1038 Yang, F., Choi, Y. J., Misch, A., Yang, X., & Dunham, Y. (2018). In defense of the commons:  
1039 Young children negatively evaluate and sanction free riders. *Psychological Science*, 29(10),  
1040 1598-1611. <https://doi.org/10.1177/0956797618779061>
- 1041 Young, L., & Saxe, R. (2011). When ignorance is no excuse: Different roles for intent across moral  
1042 domains. *Cognition*, 120(2), 202-214. <https://doi.org/10.1016/j.cognition.2011.04.005>
- 1043 Young, L., & Tsoi, L. (2013). When mental states matter, when they don't, and what that means for  
1044 morality. *Social and Personality Psychology Compass*, 7(8), 585-604.  
1045 <https://doi.org/10.1111/spc3.12044>
- 1046 Yudkin, D. A., Van Bavel, J. J., & Rhodes, M. (2020). Young children police group members at  
1047 personal cost. *Journal of Experimental Psychology: General*, 149(1), 182.  
1048 <https://doi.org/10.1037/xge0000613>
- 1049 Zelazo, P. D., Helwig, C. C., & Lau, A. (1996). Intention, act, and outcome in behavioral prediction  
1050 and moral judgment. *Child Development*, 67(5), 2478-2492. [https://doi.org/10.1111/j.1467-  
1051 8624.1996.tb01869.x](https://doi.org/10.1111/j.1467-8624.1996.tb01869.x)