



## King's Research Portal

*Document Version*  
Peer reviewed version

[Link to publication record in King's Research Portal](#)

*Citation for published version (APA):*

Barnby, J. M., Robinson, O., Deeley, P. Q., Raihani, N., Bell, V., & Mehta, M. A. (in press). Paranoia, sensitisation, and social inference: findings from two large-scale, multi-round behavioural experiments. *R Soc Open Sci.*

### **Citing this paper**

Please note that where the full-text provided on King's Research Portal is the Author Accepted Manuscript or Post-Print version this may differ from the final Published version. If citing, it is advised that you check and use the publisher's definitive version for pagination, volume/issue, and date of publication details. And where the final published version is provided on the Research Portal, if citing you are again advised to check the publisher's website for any subsequent corrections.

### **General rights**

Copyright and moral rights for the publications made accessible in the Research Portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognize and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the Research Portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the Research Portal

### **Take down policy**

If you believe that this document breaches copyright please contact [librarypure@kcl.ac.uk](mailto:librarypure@kcl.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.

# Paranoia, sensitisation, and social inference: findings from two large-scale, multi-round behavioural experiments

Barnby, J.M.<sup>1\*</sup>, Deeley, Q.<sup>2</sup>, Robinson, O.<sup>3</sup>, Raihani, N.<sup>4</sup>, Bell, V.<sup>1,5#</sup>, Mehta., M.A.<sup>1#</sup>

## Affiliations:

<sup>1</sup> Social and Cultural Neuroscience Research Group, Centre for Neuroimaging Sciences, Institute of Psychiatry, Psychology, and Neuroscience, King's College London, London, UK

<sup>2</sup> Social and Cultural Neuroscience Research Group, Forensic and Neurodevelopmental Sciences, Institute of Psychiatry, Psychology, and Neuroscience, King's College London, London, UK

<sup>3</sup> Institute of Cognitive Neuroscience, University College London, London, UK.

<sup>4</sup> Psychology and Language Sciences, University College London, London, UK.

<sup>5</sup> Research Department of Clinical, Educational, and Healthy Psychology, University College London, London, UK.

\*Corresponding Author

#These authors contributed equally to the work

## Correspondence:

Email: [joe.barnby@kcl.ac.uk](mailto:joe.barnby@kcl.ac.uk)

Twitter: @joebarnby

1 **Abstract:**

2 The sensitisation model suggests paranoia is explained by over-sensitivity to social  
3 threat. However, this has been difficult to test experimentally. We report two pre-  
4 registered social interaction studies that tested i) whether paranoia predicted overall  
5 attribution and peak attribution of harmful intent, and; ii) whether anxiety,  
6 interpersonal sensitivity, and worry predicted attribution of harmful intent. In study  
7 one, we recruited a large general population sample ( $N=987$ ) who serially interacted  
8 with other participants in multi-round Dictator games, matched to fair, partially fair, or  
9 unfair partners. Participants rated attributions of harmful intent and self-interest after  
10 each interaction. In study two ( $N=1011$ ), a new sample of participants completed the  
11 same procedure and additionally completed measures of anxiety, worry and  
12 interpersonal sensitivity. As predicted, prior paranoid ideation was associated with  
13 higher and faster overall harmful intent attributions, whereas attributions of self-  
14 interest were unaffected, supporting the sensitisation model. Contrary to predictions,  
15 neither worry, interpersonal sensitivity, nor anxiety were associated with harmful  
16 intent attributions. In a third exploratory internal meta-analysis we combined data  
17 sets to examine the effect of paranoia on trial by trial attributional changes when  
18 playing fair and unfair dictators. Paranoia was associated with a greater reduction in  
19 harmful intent attributions when playing a fair but not unfair dictator, suggesting  
20 paranoia may also exaggerate the volatility of beliefs about the harmful intent of  
21 others.

22

## 23 **1.0 Background**

24 Paranoia is a common feature in psychosis and involves an unfounded belief that  
25 others intend harm, now or in the future (1). Paranoid beliefs can be induced by  
26 recreational drugs (2,3), following sleep deprivation (4) during or after seizures (5), or  
27 from being subject to high stress (6). Paranoia also exists as a continuous trait in the  
28 general population and has shown to be characterised by interpersonal sensitivity,  
29 mistrust, ideas of reference, and ideas of persecution (7, 8).

30 Once developed, paranoid beliefs are maintained by several personal and  
31 interpersonal factors. On the personal level, worry, insomnia (9) anxiety (10, 11),  
32 probabilistic reasoning biases (12), belief inflexibility (13), and safety behaviours  
33 (avoiding the source of perceived threat) (14) all contribute to paranoia. Interpersonal  
34 cognitive biases also affect how individuals interpret social situations. The most  
35 established effect is that those with paranoid beliefs have an externalising attribution  
36 bias, whereby causes of negative events are more likely to be attributed to other  
37 people (15). Trait interpersonal sensitivity has also been associated with paranoid  
38 thinking. Those at high risk of developing psychosis report increased paranoid  
39 thinking following simulated interactions in a virtual social environment which was  
40 predicted (16) or mediated (17) by interpersonal sensitivity.

41 The sensitisation model of psychosis argues that environmental stresses and genetic  
42 vulnerabilities sensitise biological, cognitive, and affective processes to produce  
43 symptoms of psychosis, and importantly, paranoid beliefs (18, 19, 20). Neuroimaging  
44 studies have observed increased presynaptic dopamine leading up to (21) and  
45 during (22) the development of psychotic symptoms, suggesting aberrant  
46 dopaminergic transmission as crucial in sensitisation (23). Experimental data support  
47 the sensitisation of cognitive and affective processes that manifests as a 'jumping to  
48 conclusions' probabilistic reasoning bias (12, 24), high initial mistrust (25, 26) and  
49 more threatening or negatively valenced responses following heightened social  
50 arousal (27, 28).

51 One prediction arising from this model is that those high in paranoid ideation will  
52 show increased sensitivity to interpersonal interactions, and specifically potential or

53 actual social threat, leading to an increased tendency to attribute harmful intent to  
54 others, putatively both more quickly and to a greater degree.

55 Economic games derived from game theory have been previously used to test the  
56 effect of paranoia on intention attributions. These games allow for social interactions  
57 within a tightly controlled environment. Participants make decisions that have  
58 outcomes with genuine gains and losses and therefore real, albeit small, harms and  
59 benefits (29, 25). Existing research has shown that increases in harmful intent  
60 attributions are associated with trait paranoia, social threat (29, 30), social cohesion  
61 of partners in a game (31), and greater relative social rank, and outgroup status, of  
62 the interaction partner (32). However, current game theory paradigms in paranoia  
63 research that have allowed for participant-to-participant (rather than simulated; 16,  
64 17) interactions have tended to use single round games or brief interactions that are  
65 not able to test the effect of paranoia and additional psychological variables on  
66 attributions over evolving interactions.

67 In this study, we implemented a multiple-round game theory interaction using serial  
68 Dictator games. The Dictator game has been used widely in paranoia research (29,  
69 30, 32) and involves a situation where two participants are paired and one (the  
70 'dictator') is given a sum of money that they can choose to share with the 'receiver'  
71 participant (33). The receiver has no control and must accept any amount that the  
72 dictator offers. The game has been previously modified to assess social inferences  
73 made by the receiver (29, 30). After each interaction, receivers are required to rate to  
74 what extent the dictators were motivated by self-interest or an intent to harm. In the  
75 paradigm developed for this study, participants completed six serial Dictator trials  
76 against fair, partially fair and unfair partners, while rating harmful intent and self-  
77 interest motivating their partner's actions, allowing a test of sensitivity over evolving  
78 social interactions. This also allowed us to test the effect of several key affective  
79 processes previously identified as important in paranoia, namely anxiety, worry, and  
80 interpersonal sensitivity.

81 The sensitisation model of paranoia suggests several hypotheses we tested over two  
82 studies. In study one, we hypothesised that high levels of paranoid ideation would  
83 predict earlier and larger harmful intent attributions during the multi-round interaction.  
84 In study two we hypothesised that harmful intent attributions would be predicted by

85 anxiety, interpersonal sensitivity, and worry. Studies one and two were pre-  
86 registered and included hypotheses designed to replicate findings from previous  
87 studies (high attribution of harmful intent is associated with higher paranoia and  
88 unfair dictators; 29-32) as well as the key experimental hypotheses described above.  
89 Finally, we combined data from study one and two to complete exploratory analysis  
90 to gain better resolution on trial by trial effects, dictator exposure effects, and dictator  
91 behaviour overall.

92

## 93 **2.0 Study 1**

94 This study tested the main hypothesis that paranoid ideation predicts in-the-moment  
95 harmful intent attributions within serial interpersonal interactions, both in terms of  
96 overall value and by how quickly individuals reach a marker of high harmful intent  
97 attribution. Specifically, in line with prior work (29,30), we predicted that pre-existing  
98 paranoia and more unfair dictator behaviour would lead to higher harmful intent  
99 attributions. In line with prior theory (27, 28) suggesting initial sensitisation following  
100 negatively valenced responses we also predicted that paranoia would lead to fewer  
101 trials before a high harmful intent score was reached (specific preregistered  
102 predictions for this study can be found here <https://aspredicted.org/ka4ny.pdf>).

## 103 **2.1 Methodology**

104 This project was approved by the King's College London ethics board (**Study 1:**  
105 MRS-17/18-8312). All data were collected in September 2018 using Prolific  
106 Academic (hereafter Prolific; [www.prolific.ac](http://www.prolific.ac)), an online crowdsourcing platform.  
107 All data and analysis scripts are available online (<https://osf.io/u92rg/>).

108 Prior to taking part in both studies, participants were informed that their  
109 participation was voluntary, and were required to tick a box giving consent for the  
110 authors to use their anonymous data for research purposes. Using Prolific allowed  
111 rapid recruitment of a more demographically diverse sample of participants than  
112 recruitment from our social media or university networks (34). We included  
113 participants from the UK who were fluent in English and had no current or history  
114 of mental illness.

115 We recruited 987 participants (372 males). 226 people would be required to detect  
116 an effect size of 0.1 with at least 7 predictors in a multiple regression model. In  
117 order to produce robust inferences, we recruited the maximum number of  
118 participants that our resources would allow. Participants first completed the Green  
119 Paranoid Thoughts Scale (GPTS; 35). Participants were asked to indicate the  
120 extent of feelings described in 32 statements using a Likert Scale of 1 to 5, where  
121 1 = Not at All and 5 = Totally. Scores can range from 32–160, with higher scores  
122 indicating a greater degree of paranoia. The GPTS was chosen as a suitable

123 measure as it includes both core aspects of the definition of paranoia (1): social  
124 concerns about others and perception of intended harm. It has also shown to be  
125 the most reliable and valid scale for measuring paranoia across the clinical and  
126 non-clinical spectrum (36). Total paranoia scores were obtained for each  
127 participant by summing the response scores to all questions, comprising both the  
128 social reference and the persecution scales. Hereafter, this variable is referred to  
129 as 'paranoia'.

130 After completing the survey, and in keeping with Raihani and Bell (29, 30) we  
131 allowed a minimum interval of 7 days to elapse before inviting all participants to  
132 take part in the multi-round dictator game.

133 We developed a within-subjects, multi-trial modification of the Dictator game  
134 design used in previous studies to assess paranoia (See Appendix A; 29, 30).  
135 Each participant played six trials against three different types of dictator. In each  
136 trial, participants were told that they had been endowed with a total of £0.10 and  
137 their partner (the dictator) had the choice to take half (£0.05) or all (£0.10) the  
138 money from the participant. Dictators were set to either always take half of the  
139 money, have a 50:50 chance to take half or all of the money, or always take all of  
140 the money. This was noted in this study as Fair, Partially Fair, and Unfair,  
141 respectively. The order that participants were matched with dictators was  
142 randomised. Each dictator had a corresponding cartoon avatar with a neutral  
143 expression to support the perception that each of the six trials was with the same  
144 partner.

145 After each trial, participants were asked to rate on a scale of 1-100 (initialised at  
146 50) to what degree they believed that the dictator was motivated a) by a desire to  
147 earn more (self-Interest) and b) by a desire to reduce their bonus in the trial  
148 (harmful intent). Following each block of six trials, participants were asked to rate  
149 the character of the dictator overall by scoring intention again on both scales.  
150 Therefore, participants judged their perceived intention of the dictator on both a  
151 trial-by-trial and summary level.

152 After making all 42 attributions (two attributions for each of the 6 trials over 3  
153 partners, plus three additional overall attributions for each partner), participants were



154 put in the role of the dictator for 6 trials – whether to make a fair or unfair split of  
155 £0.10. Participants were first asked to choose an avatar from nine different cartoon  
156 faces before deciding on their 6 different splits. These dictator decisions were  
157 primarily collected to truthfully inform participants that decisions were made by  
158 real people (as in prior studies using this method, see 29, 30). We also included  
159 the decisions made by participants in an exploratory analysis.

160 This modification to the original dictator game design allowed us to track how  
161 changes in pre-existing paranoia were associated with changes in attributions  
162 about partner behaviour and the order of initial partner exposure, and whether  
163 attributions were highly variable over trials or consistent. 812 participants (294  
164 males) were able to be followed-up to play the multi-round dictator game. The  
165 mean age range of participants was 36-40 in the second sample.

166 All participants were paid for their completion of the GPTS, regardless of follow up.  
167 Participants were paid a baseline payment for their completion of the dictator  
168 game, along with any additional bonuses won in the game.

### 169 *Analysis*

170 Analyses conform to those outlined in our preregistration unless stated otherwise.

171 This study used an information-theoretic approach for confirmatory analysis. We  
172 analysed the data using multi-model selection with model averaging (described in  
173 29, 30). The Akaike information criterion, corrected for small sample sizes (AICc),  
174 was used to evaluate models, with lower AICc values indicating a better fit (37). The  
175 best models are those with the lowest AICc value. To adjust for the intrinsic  
176 uncertainty over which model is the true ‘best’ model, we averaged over the models  
177 in the top model set to generate model-averaged effect sizes and confidence  
178 intervals (38). In addition, parameter estimates, and confidence intervals are  
179 provided with the full global model to robustly report a variable’s effect in a model  
180 (39). This used package “MuMIn” (version 1.43.1; 40). All analyses were performed  
181 in R (version 3.6.0; 41) on an Apple OSX operating system (Mojave, 10.14.6).  
182 Visualisations were generated using the package ‘ggplot2’ (version 3.2.1; 42).

183 In our models, baseline continuous scale scores were centred and scaled to produce  
184 Z values. Model statistics reported are beta coefficients.

185 Average scores of harmful intention attributions and self-interest for each dictator  
186 were taken over each six trials for trial analysis. These were used for cumulative link  
187 mixed-models (clmm; 43). Harmful intent and self-interest attributions were set as  
188 our dependent variable. Paranoia, dictator order, dictator behaviour (fair, unfair,  
189 partially fair), age, sex, and paranoia x dictator behaviour were set as our  
190 explanatory terms with ID set as the random term.

191 For our third prediction, participants that scored above 60 were considered to have  
192 “high” harmful intent attributions. Both harmful intent and self-interest scores  
193 participants were set a value of 6 if they had scored 60 in their first trial, 5 if they had  
194 scored over 60 by their second trial, 4 if they had scored 60 by their third trial, and so  
195 on. We report this result, but also wanted to consider a high harm attribution as  
196 someone that scored over the mean harmful intent attribution relative to every other  
197 participant in the population for each dictator. This latter analysis reported in addition  
198 to our preregistered plan, which was based on previous mean group estimates.  
199 Mean thresholds for each dictator are stated for each analysis in the Results. All  
200 trials following the threshold being reached were coded as 0. Participants not  
201 reaching the threshold for any trial were coded 0 across all trials. Both unfair and fair  
202 dictator behaviour were analysed with two cumulative link models (clm) each, one for  
203 harm-intent and one for self-interest. This slightly deviates from our preregistration  
204 that suggests the use of Kruskal-Wallis and Dunn post-hoc tests, however we  
205 decided that using a clm was a more robust way to analyse the data.

206 For visualisation purposes we calculated paranoia groups based on the quantiles of  
207 GPTS scores across the population, and additionally divided those in the top quantile  
208 by those exceeding the clinical mean of paranoia defined in previous work (101.9;  
209 35). These divisions were: Low (<36; n = 232) Medium (36-43; n = 180), High (44-  
210 59; n = 199), and Very High (59-101.9; n = 167), and Clinical (>102, n = 34). This  
211 variable is hereafter named paranoia ‘level’. Slightly different score parameters for  
212 each paranoia level were included in our pre-registration but we have adapted them  
213 in this study based on our population GPTS quartiles, although GPTS is maintained  
214 as a continuous term in our models.

215

## 216 **2.2 Results**

217 812 participants that were able to be followed up were included in the analysis. 15  
218 were removed for incomplete data, 24 removed for failing both control questions,  
219 and 136 for non-participation in the multi-round dictator game. Mean baseline  
220 paranoid ideation in the excluded participants ( $M = 50.43$ ,  $SE = 1.62$ , range = 32-  
221 134) were comparable to those that were included in the analysis ( $t(252) = 0.322$ ,  
222 95%CI: -2.93, 4.08).

### 223 *Explanatory variables of baseline paranoia score*

224 Paranoia scores ranged from 32.0-149.0 with a mean of 51.0 (SE: 0.74; Skew: 1.7).  
225 Older participants were less paranoid (-1.89; 95% CI: -2.22, -1.57), male participants  
226 were more paranoid (0.17; 95% CI: 0.04, 0.34), and there was no effect of education  
227 on paranoia (-0.39; 95% CI: -1.16, 0.17).

### 228 *Prediction 1: Paranoia and harmful intent*

229 As predicted, paranoia positively predicted higher HI attributions across all three  
230 dictators, however, there was no effect of paranoia on SI attributions (Table 1).

### 231 *Prediction 2: Dictator behaviour and harmful intent*

232 As predicted, as dictators were increasingly unfair (higher proportion of unfair  
233 decisions), higher HI and SI attributions were observed (Table 1). Figure 1a depicts  
234 the difference in HI and SI attributions between the population when delineated by  
235 their paranoia level (low, medium, high, very high, clinical) for Study 1.

### 236 *Prediction 3: Paranoia and earlier high harmful intent attributions*

237 As predicted, high (over 60) harmful intent attributions were triggered in earlier trials  
238 as paranoia increases for both unfair (-0.12; 95% CI: -0.21, -0.03) and fair (-0.14,  
239 95% CI: -0.33, -0.01) dictators, however this was not found for SI attributions (see  
240 Appendix B).

### 241 *Exploratory analysis*

242 We also completed an analysis using a relative threshold for earlier high decisions  
243 based on the mean of the population for each dictator rather than a pre-set cut-off of

244 60 as in the preregistered analysis. For unfair dictators, high (mean = 53.51) HI  
245 attributions were triggered in earlier trials as paranoia increased (-0.12; 95% CI: -  
246 0.20, -0.02). However, this was not found for fair dictators (mean = 24.26) (-0.06;  
247 95% CI: -0.19, 0.01). This was not found for SI attributions in either dictator  
248 condition. See figure 2a for trial-by-trial average attributions across participants for  
249 study 1.

250

251 **Table 1. Variables affecting Harmful Intention and Self Interest scores in the**  
 252 **multi-round dictator game (Study 1).** Harmful Intent was coded as a five-level  
 253 ordinal categorical variable and set as the response term in the clmm. Participant ID  
 254 was set as the random variable (43). Relative Importance is the probability that the  
 255 term in question is a component of the true best model and a value for the amount of  
 256 times the term is included in the selection of top models to be averaged. Order refers  
 257 to the order in which a fair, partially fair, or unfair dictator was presented to  
 258 participants. An interaction between dictator and paranoia is not included in the  
 259 model for Self Interest Attributions as it was not included in the final top model. Age  
 260 was not included in the final top model and is therefore absent from the tables.

Parameter	Estimate	Standard Error	95% CI		Relative Importance
			Lower	Upper	
<b>Harmful Intent Attributions</b>					
<i>Intercept 1 2</i>	-1.26	0.11	-1.48	-1.05	
<i>Intercept 2 3</i>	0.47	0.10	0.27	0.68	
<i>Intercept 3 4</i>	2.17	0.12	1.94	2.39	
<i>Intercept 4 5</i>	3.67	0.14	3.41	3.94	
<b>Dictator</b> (Fair < Partially Fair < Unfair)	2.22	0.09	2.06	2.39	1
<b>Order</b> (Fair < Partially Fair < Unfair)	-1.12	0.15	-1.42	-0.83	1
<b>Paranoia (Z score)</b>	0.36	0.09	0.19	0.53	1
<b>Sex (Male   Female)</b>	-0.03	0.11	-0.26	0.19	0.25
<b>Dictator x Paranoia</b>	0.14	0.10	-0.06	0.34	0.79
<b>Self Interest Attributions</b>					
<i>Intercept 1 2</i>	-6.53	0.25	-7.01	-6.05	
<i>Intercept 2 3</i>	-5.25	0.21	-5.66	-4.84	
<i>Intercept 3 4</i>	-3.15	0.16	-3.46	-2.84	
<i>Intercept 4 5</i>	-0.28	0.11	-0.50	-0.07	
<b>Dictator</b> (Fair < Partially Fair < Unfair)	4.33	0.17	3.99	4.67	1
<b>Order</b> (Fair < Partially Fair < Unfair)	-0.82	0.16	-1.13	-0.50	1
<b>Paranoia (Z score)</b>	0.01	0.05	-0.09	0.11	0.24
<b>Sex (Male   Female)</b>	-0.03	0.11	-0.23	0.18	0.23

261

262

## 263 **3.0 Study 2**

264 Study 1 suggested that prior paranoid beliefs led to larger and earlier harmful intent  
265 attributions. Prior models and evidence (9, 10, 11, 16, 17) suggest a role of affective  
266 processes in paranoid ideation, specifically interpersonal sensitivity (44), state and  
267 trait anxiety (45) and worry (46). Therefore, we also wanted to test whether these  
268 psychological variables predict harmful intent attributions. Specifically, we predicted  
269 that scores on measures of state anxiety and overall interpersonal sensitivity (and  
270 specifically subscales of “Fragile Inner Self” and “Interpersonal Awareness”) would be  
271 associated with higher harmful intent attributions. We also predicted that state  
272 anxiety and pre-existing paranoia would interact and be associated with higher and  
273 earlier scores of harmful intent attributions (specific preregistered predictions for this  
274 study can be found here <http://aspredicted.org/yz5gr.pdf>).

## 275 **3.1 Methodology**

276 This project was approved by the King’s College London ethics board (Study 2: LRS-  
277 18/19-9281). Data were collected in February 2019 using Prolific. Data and analysis  
278 scripts are available online (<https://osf.io/u92rg/>).

279 We recruited 1011 participants (374 males). 226 people would be required to detect  
280 an effect size of 0.1 with at least 7 predictors in a multiple regression model. In  
281 order to produce robust inferences, we recruited the maximum number of  
282 participants that our resources would allow. Participants recruited for this study  
283 were not participants in Study 1. Study procedures and analyses were identical to  
284 Study 1 aside from the inclusion of anxiety, worry and interpersonal sensitivity  
285 measures.

286 We assessed both trait anxiety and state anxiety using the State-Trait Anxiety  
287 Inventory (STAI; 45). It is comprised of two subscales, one for trait and one for state  
288 anxiety, each made of 20 items. Each item is rated on a scale of one to four, from  
289 “Almost Never” to “Almost Always”. The trait measure was given to participants at  
290 baseline alongside the GPTS. The state measure was given immediately after the  
291 multi-round dictator game.

292 We measured interpersonal sensitivity using the Interpersonal Sensitivity Measure  
293 (ISM; 44). The ISM is comprised of five subscales: Fragile Inner Self (5 items), Need  
294 for Attachment (8 items), Interpersonal Awareness (7 items), Timidity (8 items), and  
295 Separation Anxiety (8 items). Each item is on a scale of one to four, from “Very  
296 Unlike You” to “Very Like You”. Subscales are summed to form summary scores.  
297 The ISM was given at baseline alongside the GPTS.

298 We also measured worry using the Penn-State Worry Questionnaire (PSWQ) (46) as  
299 worry has been additionally implicated as highly predictive of paranoia (1). The  
300 PSWQ is comprised of 16 items, each on a scale of one to five, from “Not at all  
301 typical of me” to “Very typical of me”. The PSWQ was given at baseline alongside  
302 the GPTS.

303 All analyses were performed in R (version 3.6.0; 41) on an Apple OSX operating  
304 system (Mojave, 10.14.6).

305 Analyses conform to our preregistration unless stated otherwise. We included the  
306 explanatory variables from the STAI, PSWQ and ISM in our cumulative link mixed  
307 models alongside the GPTS scores with ID set as the random variable. Continuous  
308 variables were z-score transformed. Model statistics reported are beta coefficients  
309 unless stated otherwise.

310 Notably, we ran extra exploratory analyses (not preregistered) to assess the  
311 association between extra explanatory variables (STAI, PSWQ, and ISM) and  
312 attributions when paranoia was excluded from the model to account for any potential  
313 suppressing effects of paranoia in the models that may exaggerate associations  
314 between anxiety, worry, or interpersonal sensitivity and attributions.

315

316 **3.2 Results**

317 885 participants that were able to be followed up were included in the analysis. 8  
318 were removed for incomplete data and 118 for non-participation in the multi-round  
319 dictator game. Mean baseline paranoid ideation in the excluded participants ( $M =$   
320  $58.54$ ,  $SE = 2.35$ , range = 32-140) were higher than participants that were  
321 included in the analysis ( $t(153) = -2.41$ , 95%CI: -10.85, -1.09) by a small amount.

322 *Explanatory variables of baseline paranoia*

323 Paranoia scores ranged from 32-159 with a mean of 53 (SE: 0.45; Skew: 1.54).  
324 Older participants were less paranoid (-0.05; 95% CI: -0.05, -0.04), there was a  
325 negligible effect of being male on paranoia (0.05; 95% CI: -0.04, 0.24), and there  
326 was a quadratic (-1.20, 95%CI: -1.80, -0.60) relationship between education and  
327 paranoia. Paranoia positively correlated with anxiety, worry, and interpersonal  
328 sensitivity ( $R = 0.38-0.51$ , see Appendix C).

329 *Replication of main findings of study 1*

330 Paranoia positively predicted higher HI attributions across all three dictators, there  
331 was no effect of paranoia on SI attributions, and additionally, unfairness of dictator  
332 was associated with higher HI and SI attributions. Order effects were also replicated  
333 (See Figure 1 and Table 2).

334 For unfair dictators, high (mean = 46.56) HI attributions were not uniformly observed  
335 in earlier trials as paranoia increased (-0.06; 95% CI: -0.17, 0.01), but were for fair  
336 dictators (mean = 21.39) (-0.12; 95% CI: -0.20, -0.03). However, paranoia was not  
337 associated with high SI attributions in earlier trials in either dictator condition.

338 Figure 2b shows average trial-by-trial attributions for each level of paranoia in **Study**  
339 **2**.

340



341 **Table 2. Variables affecting Harmful Intent and Self Interest scores in the multi-**  
 342 **round dictator game (Study 2).** Harmful Intent was coded as a five-level ordinal  
 343 categorical variable and set as the response term in the clmm. Participant ID was set  
 344 as the random variable (43). Relative Importance is the probability that the term in  
 345 question is a component of the true best model and a value for the amount of times  
 346 the term is included in the selection of top models to be averaged. Order refers to the  
 347 order in which a fair, partially fair, or unfair dictator was presented to participants.

Parameter	Estimate	Standard Error	95% CI		Relative Importance
			Lower	Upper	
<b>Harmful Intent Attributions</b>					
<i>Intercept 1 2</i>	-0.64	0.23	-1.09	-0.18	
<i>Intercept 2 3</i>	1.28	0.24	0.82	1.74	
<i>Intercept 3 4</i>	2.95	0.25	2.47	3.43	
<i>Intercept 4 5</i>	4.38	0.26	3.88	4.89	
<b>Dictator</b> (Fair < Partially Fair < Unfair)	2.00	0.09	1.82	2.18	1
<b>Order</b> (Fair < Partially Fair < Unfair)	-1.17	0.17	-1.52	-0.83	1
<b>Paranoia</b> (Z score)	0.35	0.10	0.15	0.54	1
<b>Sex</b> (Male   Female)	-0.16	0.21	-0.71	0.10	0.52
<b>Age</b>	0.00	0.01	-0.01	0.02	0.32
<b>Self Interest Attributions</b>					
<i>Intercept 1 2</i>	-6.59	0.35	-7.27	-5.91	
<i>Intercept 2 3</i>	-5.35	0.33	-5.99	-4.71	
<i>Intercept 3 4</i>	-3.16	0.30	-3.75	-2.58	
<i>Intercept 4 5</i>	-0.21	0.28	-0.75	0.33	
<b>Dictator</b> (Fair < Partially Fair < Unfair)	4.59	0.17	4.26	4.93	1
<b>Order</b> (Fair < Partially Fair < Unfair)	-0.71	0.16	-1.02	-0.39	1
<b>Paranoia</b> (Z score)	-0.03	0.07	-0.28	0.09	0.34
<b>Sex</b> (Male   Female)	0.01	0.07	-0.31	0.43	0.11
<b>Age</b>	0.00	0.01	-0.02	0.00	0.43

348  
349

350

[FIGURE 1 HERE]

351

352 **Figure 1. Average Self-Interest attributions (Blue) and Harmful Intent**  
353 **attributions (Red), averaged across trials for divisions of GPTS score and**  
354 **faceted by each type of dictator for both study 1 (A;  $n = 812$ ) and study 2 (B;  $n$**   
355 **= 885). Dots represent the mean for each level of paranoia. Lines represent the**  
356 **95% confidence interval. Participants played against different partners in a pseudo-**  
357 **random order. 'Clinical' refers to participants in the general population who scored**  
358 **past the threshold for GPTS scores typical in clinical populations (101.9; 35).**

359

360

361

[FIGURE 2 HERE]

362

363 **Figure 2. Average Harmful Intent (Red) and Average Self-Interest (Blue)**  
364 **attributions for each trial across divisions of GPTS scores, faceted by type of**  
365 **dictator for study 1 (A;  $n = 812$ ) and study 2 (B;  $n = 885$ ).** Points = mean, bars =  
366 95% confidence interval. Grey lines = mean score across the group. 'Clinical' refers  
367 to participants in the general population who scored past the threshold for GPTS  
368 scores typical in clinical populations (101.9; 35).

369

370 *Predictions 1 and 2: State anxiety, paranoia and harmful intent*

371 Contrary to predictions, state anxiety did not predict overall HI or SI attributions in  
372 any dictator condition and there was no interaction with paranoia. Excluding paranoia  
373 did not change the conclusions of the estimates (Table 3).

374 *Prediction 3 and 4: Interpersonal sensitivity, paranoia and harmful intent*

375 Contrary to predictions, interpersonal sensitivity predicted a decrease in overall HI (-  
376 0.29, 95%CI: -0.49, -0.10) but not SI attributions across all dictators, and there was  
377 no interaction between interpersonal sensitivity and paranoia for HI or SI attributions  
378 across all dictators. However, excluding paranoia from the models removed the  
379 association of interpersonal sensitivity on HI attributions (Table 3), indicating a  
380 suppressing effect of paranoia.

381 We ran a model with all subscales of interpersonal sensitivity included. This  
382 suggested that 'interpersonal awareness' (-0.54, 95%CI: -0.80, -0.28) was negatively  
383 associated with HI attributions while 'separation anxiety' (0.36, 95%CI: 0.08, 0.64)  
384 was positively associated. Conversely, 'timidity' was negatively associated with SI  
385 attributions (-0.46, 95%CI: -0.68, -0.23) and 'interpersonal awareness' (0.31, 95%CI:  
386 0.04, 0.58) and 'need for attachment' (0.28, 95%CI: 0.07, 0.48) were positively  
387 associated. Full Model statistics that included all predictors together with paranoia  
388 can be found in Appendix D.

389 We also ran all subscales of the interpersonal sensitivity measure in separate  
390 models. The "need for attachment" subscale of the ISM was associated with a  
391 decrease in HI attributions, and all other subscales had no effect. 'Timidity' was  
392 associated with reduced SI attributions, and 'need for attachment' was associated  
393 with an increase in SI scores (Table 3).

394 *Prediction 5: Anxiety, paranoia, and trials to peak decision*

395 Contrary to predictions, state anxiety alone and its interaction with paranoia did not  
396 predict scoring above the mean in an earlier trial for HI and SI attributions during  
397 both unfair and fair dictators. Excluding paranoia from the models did not change the  
398 conclusions of the estimates (Table 3).

399 Paranoia also remained an independent predictor of harmful intent regardless of  
 400 anxiety, interpersonal sensitivity, or worry being included in models (0.34-0.60,  
 401 95%CI: 0.13-0.58, 0.54-0.88; see Appendix D).

402 **Table 3. Summary of extra explanatory variables affecting Harmful Intention**  
 403 **and Self Interest attributions in the multi-round dictator game (Study 2).**  
 404 Harmful Intent was coded as a five-level ordinal categorical variable and set as the  
 405 response term in the clmm. Participant ID was set as the random variable (42). All  
 406 predictors were run in separate models with Dictator, Age, and Sex as other fixed  
 407 effects. NA = not included in the final top model. Appendix D contains all the  
 408 estimates of predictors in the models that included paranoia.

409

Parameter	Estimate	Standard Error	95% CI		Relative Importance
			Lower	Upper	
<b>Harmful Intent Attributions</b>					
Trait Anxiety	0.01	0.03	-0.04	0.07	0.19
State Anxiety	0.04	0.06	-0.09	0.18	0.37
Interpersonal Sensitivity	-0.01	0.02	-0.06	0.04	0.19
Interpersonal Awareness	-0.13	0.11	-0.36	0.09	0.70
Separation Anxiety	0.08	0.14	-0.13	0.28	0.51
Timidity	-0.11	0.12	-0.34	0.12	0.63
Need for Attachment	-0.22	0.10	-0.42	-0.01	1
Fragile Inner Self	0.07	0.09	-0.13	0.26	0.49
Worry	0.06	0.10	-0.13	0.24	0.42
<b>Self Interest Attributions</b>					
Trait Anxiety	0.01	0.05	-0.11	0.25	0.17
State Anxiety	0.09	0.1	-0.12	0.29	0.57
Interpersonal Sensitivity	NA	NA	NA	NA	NA
Interpersonal Awareness	0.03	0.08	-0.11	0.17	0.34
Separation Anxiety	-0.00	0.05	-0.10	0.09	0.26
Timidity	-0.20	0.09	-0.38	-0.02	1
Need for Attachment	0.20	0.09	0.02	0.38	1
Fragile Inner Self	-0.03	0.06	-0.15	0.10	0.35
Worry	0.11	0.11	-0.11	0.32	0.66

410

[FIGURE 3 HERE]

411

412 **Figure 3. Pearson R correlations for centred and scaled scores on state and**  
413 **trait anxiety, paranoia, interpersonal sensitivity, and worry questionnaires by**  
414 **harmful intent (A) and self-interest (B) scores in Study 2, faceted by dictator**  
415 **condition.  $N = 885$ .**

## 416 **4.0 Internal Meta-Analysis**

417 We combined data from Study 1 and 2 to analyse the overall effect of paranoia, trial  
418 by trial attributional change for each dictator, as well as order effects, and overall  
419 dictator behaviour on attributions. We also include an exploratory analysis  
420 recommended by a reviewer to assess whether overall harmful intent and self-  
421 interest attributions made across partners in the task, and pre-existing paranoia,  
422 affected unfair decisions made by participants when they took the role of the dictator  
423 with a new partner.

### 424 **4.1 Methodology**

425 A total of 1754 participants were included in the analysis from **Study 1** and **Study 2**.  
426 The meta-analysis was not preregistered, although data and analysis scripts are  
427 available online (<https://osf.io/u92rg/>).

428 As in both previous studies, paranoia scores on the GPTS were divided into  
429 quantiles (Low, 32-36; medium, 37-44; high, 45-61; very high, 61 – 101.9) and also a  
430 group who passed GPTS scores exceeding the clinical mean (clinical, >101.9) (See  
431 Figure 4).

432 All analyses were performed in R (version 3.6.0; 41) on an Apple OSX operating  
433 system (Mojave, 10.14.6).

434 Linear mixed effects models (function “lmer”; package “lme4”; 47, ID as the random  
435 variable) were run to determine the effect of initial dictator exposure on overall HI  
436 and SI attributions for fair and unfair dictators. They were also used to calculate  
437 changes in HI and SI attributions for each trial relative to the first, and the overall  
438 effect of paranoia and sex on attributions. Probability distributions and uncertainty  
439 estimates of the direction of beta coefficients produced by mixed effect models were  
440 computed for HI and SI attributions for each trial and each dictator (using “rstanarm”,  
441 ID set at the random variable; version 2.18.2; 48; probability of direction fitted with  
442 “bayestestR”; version 0.3.0; 49) to give a visual description of changes in HI and SI  
443 scores as trials continued (figure 5).

444 We calculated the trial where a high (> mean) attribution was made and trial by trial  
445 changes to attributions when considering pre-existing paranoia (GPTS score).  
446 Cumulative link models with multimodal averaging (as with Study 1 and 2) were used  
447 for each dictator. Trial by trial analyses between levels of paranoia were visualised  
448 separately for harmful intent and self-interest attributions for each dictator (Figure 6).

449 Finally, we ran an exploratory analysis on the combined datasets to establish  
450 whether pre-existing paranoia, overall harmful intent attributions, overall self-interest  
451 attributions, and sex was associated with more unfair decisions made by participants  
452 when they took the role of the dictator following being the receiver. It was clear in the  
453 task that participants were making decisions for a new partner, as opposed to the  
454 partners they had been paired with in the previous trials. Participants made six  
455 dictator decisions in total. We used a mixed effects binomial regression model (using  
456 package “lme4”; version 1.1.21; 47) to assess this question, with ID and Decision  
457 Trial (1-6) as random effects. The model was unable to converge with both overall  
458 harmful intent and self-interest attributions included, so we ran separate models that  
459 included each.

460



461

[FIGURE 4 HERE]

462

463 **Figure 4. Rainbow cloud plot for each quartile of the Green Paranoid Thoughts**  
464 **Scale (GPTS).** The highest quartile was subdivided into those who had and had not  
465 passed the clinical threshold (101.9) (35). The clinical division is denoted by a grey  
466 line.

467 **4.2 Results**

468 See Appendix E for the density distributions of scores for each dictator and trial.

469 *Order effects*

470 Being initially exposed to a more unfair dictator predicted a decrease in HI  
471 attributions for fair (-3.61, 95%CI: -4.38, -2.85) and unfair dictator conditions (-16.70,  
472 95%CI: -19.50, -13.84) in the context of the whole population. Being initially exposed  
473 to a more unfair dictator predicted a decrease in self-interest attributions when  
474 playing fair (-5.89, 95%CI: -8.05, -3.74) and unfair dictator conditions (-1.66, 95%CI:  
475 -2.61, -0.71). Paranoia predicted an increase in HI attributions for both dictators in  
476 these models (fair dictator: 1.92, 95%CI: 0.91, 2.94; unfair dictator: 3.47, 95%CI:  
477 1.84, 5.11), but not SI attributions.

478 *Trial by trial analysis*

479 See Figure 5 (Appendix F for confidence intervals) for overall changes in HI and SI  
480 scores for each dictator from trials 1-6 across the population.

481 Paranoia predicted earlier trials in which a high HI score (> mean) was triggered for  
482 both unfair (-0.08, 95%CI: -0.14, -0.01) and fair (-0.08, 95%CI: -0.14, -0.02) dictators,  
483 although this was not true for SI scores. Additionally, paranoia predicted an overall  
484 decrease in scores between the first and the sixth trial for fair (-0.70, 95%CI: -1.54, -  
485 0.03) but not unfair dictators, and this was not true for SI scores for either dictator  
486 (Figure 6 for visual summary).

487 *Dictator Decision Analysis*

488 Paranoia predicted more unfair decisions made by participants (0.66, 95%CI: 0.31,  
489 1.02), as did being male (2.19, 95%CI: 1.44, 2.95) and overall self-interest  
490 attributions (0.63, 95%CI: 0.31, 0.96), whereas overall harmful intent attributions  
491 (0.27, 95%CI: -0.07, 0.62) and the order in which participants had been partnered  
492 with dictators (0.035, 95%CI: -0.25, 0.96) did not affect decisions.

493

494

[FIGURE 5 HERE]

495

496 **Figure 5: Probability distributions of beta coefficient from linear mixed effects**  
497 **models representing HI and SI attributions of the whole population by unfair**  
498 **and fair dictators between trials two to six when compared with trial one.**

499 Probability distributions of beta coefficients modulated by paranoia (zPara; scaled  
500 and centred GPTS scores) and being a male (SexMale) when compared with being  
501 a female are also included.

502

[FIGURE 6 HERE]

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

**Figure 6: Each plot displays mean and SD for harmful intent (red) and self-interest (blue) attributions, faceted by dictator, graded in colour by paranoia division.** Group comparison significant values represent HI/SI score ~ Paranoia for each trial. \* =  $p < 0.05$ , \*\* =  $p < 0.01$ , ns = not significant.

## 519 **4.0 Discussion**

520 We undertook two studies to test the sensitisation model of paranoia using a multi-  
521 round Dictator game. This controlled experimental design models social inferences  
522 about the intentions of a 'dictator' (playing partner) over successive interactions and  
523 varying conditions of fair behaviours. In study one we tested the effect of self-  
524 reported paranoid beliefs on the attribution of harmful intent. In study two we tested  
525 the effect of anxiety, worry and interpersonal sensitivity in modulating these effects.

526 In line with our predictions, paranoia was associated with earlier and higher levels of  
527 harmful intent attribution across all conditions, and higher levels of harmful intent  
528 attribution as partners were increasingly unfair in their division of resources. Contrary  
529 to predictions, we found no meaningful effects of anxiety, interpersonal sensitivity or  
530 worry on the attribution of harmful intent. Self-interest attributions were only  
531 modulated by dictator type and the order in which participants were partnered with  
532 dictators. An internal meta-analysis highlighted that paranoia was associated with  
533 greater reductions of harmful intent attributions in fair dictator conditions over six  
534 trials, but not unfair dictators, where harmful intent attributions remained consistent.  
535 Additionally, in the meta-analysis regardless of paranoia, harmful intent attributions  
536 increased over trials with unfair dictators and decreased over trials in fair dictators.  
537 Finally, prior paranoia, but not in-the-moment attributions of harmful intent, predicted  
538 more selfish decisions with a new partner in line with prior evidence (30).

539 Our data provides additional evidence for the sensitisation model in paranoia. Our  
540 findings converge with previous game theory studies on paranoia that measured  
541 attribution of harmful intent using between-subject single shot designs. In previous  
542 studies that used Dictator games, paranoia predicted greater harmful intent  
543 attributions regardless of partner fairness (29, 30). This new study replicated these  
544 findings and additionally showed through the use of a within-group design and serial  
545 interactions that paranoia was associated with faster and larger attributions of  
546 harmful intent relative to partner fairness, suggesting increased sensitivity to  
547 perceived threat in interpersonal interactions. This is in line with previous findings  
548 from studies using a range of alternative paradigms. Simulated social exclusion with  
549 the 'cyberball' game increased state paranoia in non-clinical individuals with high trait  
550 paranoia (50), in individuals at high risk of psychosis (51), and patients with paranoid

551 delusions (52). Experience sampling studies have found that moments of subjective  
552 stress (53, 54, 56) and physiological arousal (55) predict an increase in paranoia.  
553 Similarly, immersion in a stressful social environment, either in virtual reality (56) or a  
554 genuine city street (57), increased state paranoia. Additionally, our results of overall  
555 higher self-interest attributions at all levels of paranoia that are only modulated by  
556 dictator type are consistent with prior evidence (29-31), demonstrating the specificity  
557 of prior paranoia on momentary inferences instead of a general bias in social  
558 reasoning.

559 Our data also converge with theories of social learning. Models of social impression  
560 formation in healthy populations suggest that impressions of 'bad' others are more  
561 volatile, and hence updated more quickly when a putatively bad agent becomes  
562 fairer (58). Our findings that paranoia was associated with higher initial baseline  
563 harmful intent attributions, and also greater reductions in harmful intent attributions in  
564 fair partner conditions, provides convergent evidence that pre-existing paranoia may  
565 both lead to higher baseline impressions of harmful intent and concurrently amplify  
566 belief volatility.

567 Counter to our predictions, we did not find any effect of anxiety or worry on the  
568 attribution of harmful intent. Cognitive models of paranoia (59-61) cite worry and  
569 anxiety as maintaining paranoid ideation based on a range of prior evidence. Worry  
570 has been found to be present at high levels in highly paranoid people (62) and  
571 psychological treatment for worry has been shown to reduce paranoia in a targeted  
572 randomised controlled trial (63). Similarly, induction of stress has been shown to  
573 increase state paranoia, mediated by anxiety (6, 57), in addition to anxiety predicting  
574 higher state paranoia in ambiguous virtual environments (64). Given the strength of  
575 prior evidence we think it unlikely that anxiety and worry play no part in paranoia and  
576 suggest two possibilities for why no effect was found in this study. The first may be  
577 that we measured harmful intent attributions for specific events and general worry  
578 and anxiety may be more involved in maintaining paranoid ideation (i.e. promoting  
579 paranoid rumination) than amplifying in-the-moment paranoid attributions. Indeed,  
580 currently models of paranoia suggest anxiety and worry are maintenance factors for  
581 paranoid thoughts (17) and worry and experience sampling studies suggest that  
582 proximal worry and anxious rumination have a larger effect on paranoia than in time-

583 lagged analyses (65). Thus, in-the-moment attributions of harmful intent may be  
584 more dependent on momentary worry and anxiety prior to social interactions, and it  
585 is not certain that traits will become relevant to live social inferences. Secondly, other  
586 predisposing factors (e.g. trauma; 17) not measured may be more relevant to the  
587 relationship between general anxiety and harmful intent attributions.

588 Contrary to our prediction we found that interpersonal sensitivity was not associated  
589 with harmful intent attributions. A recent systematic review reported a strong  
590 relationship between interpersonal sensitivity and trait paranoia, but a variable and  
591 unclear relationship with state paranoia (66). For example, using a general  
592 population sample, virtual reality studies have found an association between state  
593 paranoia and overall interpersonal sensitivity (67), even when adjusting for  
594 confounders (68, 69). However, when using 'real world' stooges, an association with  
595 state anxiety was only found with the separation anxiety subscale (70). However, we  
596 did not find a positive relationship between harmful intent attributions and the  
597 interpersonal sensitivity measure when we included it alone without paranoia,  
598 although there was a negative association between total interpersonal sensitivity and  
599 harmful intent when paranoia was included in models. Like anxiety and worry, it may  
600 be that the influence of trait interpersonal sensitivity on momentary paranoia is  
601 dependent on different immediate social circumstance e.g. when in the presence of  
602 another person (67-69), or alternatively that interpersonal sensitivity may only relate  
603 to maintaining paranoid thoughts and not momentary harmful intent attributions.

604 We also note some limitations to this study. As with previous designs, our study used  
605 crowd-sourcing platforms. This affords us a much larger, more representative  
606 sample than university or community samples (34), with higher response rates (71),  
607 greater experimental naivety and larger chances of replication (72), although our  
608 data drew solely on a UK population. However, given our exclusion criterion  
609 (participants had to fail both control questions to be removed), it is possible that  
610 some participants did not respond accurately due to poor attention, potentially  
611 leading to inflated effect sizes (34). We note however that previous studies have  
612 found online participants to produce equal or better-quality data than lab participants  
613 for the same task (73). Additionally, it is not clear to what extent those who score  
614 above the clinical mean on the paranoia scale in this study resemble patients with

615 paranoid delusions. Given such a large sample, it would be surprising if at least  
616 some of the high scorers did not have delusions, although it is also the case that  
617 those most disabled by psychosis may be least able to participate in computer-based  
618 studies.

619 Our game theory paradigm measured harmful attributions in ambiguously motivated,  
620 loss-inducing, online interaction. One potential limitation is the extent to which  
621 participants were sceptical and believed they were being deceived by the  
622 experimenters. We found no relationship between scepticism and harmful intent  
623 attributions, and likewise our findings have replicated previous evidence using a  
624 similar manipulation (29, 30). One additional question is the extent to which our  
625 findings generalise to diverse social situations. As noted above, the results reported  
626 here converge with those reported in experience sampling studies of everyday  
627 interactions and immersive experimental studies, suggesting they also reflect the  
628 operation of common cognitive mechanisms. However, the specific differences in  
629 how paranoia manifests in online and offline contexts has yet to be tested and we  
630 feel this is something that needs further research.

## 631 **5.0 Conclusion**

632 We have demonstrated that paranoid ideation leads to quicker and exaggerated  
633 attributions of harmful intent, but not attributions of self-interest, in a motivationally  
634 ambiguous, live online social task. Our findings support the theory of sensitisation in  
635 paranoia - specifically, that pre-existing paranoid beliefs reflect a heightened  
636 sensitivity to social stress which increases attributions of harmful intent. We also  
637 show in a within group design that the cognitive processes involved in detection of  
638 social threat through fairness are at least partially distinct. The finding that anxiety,  
639 interpersonal sensitivity overall and worry did not predict attributions of harmful intent  
640 suggests that general anxiety, interpersonal sensitivity as a single measure and  
641 worry may mediate paranoid rumination rather than in-the-moment attributions.  
642 Future studies will employ game theory paradigms in patient groups to investigate  
643 the relationship between clinical and non-clinical paranoia. At a neural level,  
644 evidence of the involvement of the mesolimbic dopamine system in psychosis  
645 suggest that future studies should investigate how dopamine modulates threat  
646 attribution in illness and health.



## **Acknowledgments**

We would like to thank Uri Hertz for kindly sending his avatar images for use in this game.

## **Conflict of Interest Statement**

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

## **Funding**

JMB is supported by the UK Medical Research Council (MR/N013700/1) and King's College London member of the MRC Doctoral Training Partnership in Biomedical Sciences.

## **Author Contributions**

JMB initially devised the studies. JMB constructed the multi-round dictator game. JMB and NR revised the multi-round dictator game. JMB collected the data, analysed the data and wrote initial the draft of the manuscript. JMB, QD, OR, NR, VB and MAM critically revised the manuscript.

## 6.0 References

1. Freeman, D., & Garety, P. A. (2000). Comments on the content of persecutory delusions: Does the definition need clarification? *British Journal of Clinical Psychology*. <https://doi.org/10.1348/014466500163400>
2. Englund, A., Morrison, P. D., Nottage, J., Hague, D., Kane, F., Bonaccorso, S., ... & Feilding, A. (2013). Cannabidiol inhibits THC-elicited paranoid symptoms and hippocampal-dependent memory impairment. *Journal of Psychopharmacology*, 27(1), 19-27.
3. McKetin, R. (2018). Methamphetamine psychosis: insights from the past. *Addiction*, 113(8), 1522-1527.
4. Reeve, S., Emsley, R., Sheaves, B., & Freeman, D. (2017). Disrupting sleep: the effects of sleep loss on psychotic experiences tested in an experimental study with mediation analysis. *Schizophrenia bulletin*, 44(3), 662-671.
5. Elliott, B., Joyce, E., & Shorvon, S. (2009). Delusions, illusions and hallucinations in epilepsy: 2. Complex phenomena and psychosis. *Epilepsy research*, 85(2-3), 172-186.
6. Lincoln, T. M., Peter, N., Schäfer, M., & Moritz, S. (2009). Impact of stress on paranoia: an experimental investigation of moderators and mediators. *Psychological medicine*, 39(7), 1129-1139.
7. Bell, V., & O'Driscoll, C. (2018). The network structure of paranoia in the general population. *Social psychiatry and psychiatric epidemiology*, 53(7), 737-744.
8. Bebbington, P. E., McBride, O., Steel, C., Kuipers, E., Radovanovič, M., Brugha, T., ... Freeman, D. (2013). The structure of paranoia in the general population. *British Journal of Psychiatry*, 202(6), 419-427. <https://doi.org/10.1192/bjp.bp.112.119032>
9. Freeman, D. (2007). Suspicious minds: The psychology of persecutory delusions. *Clinical Psychology Review*, 27(4), 425-457. <https://doi.org/10.1016/j.cpr.2006.10.004>
10. Startup, H., Freeman, D., & Garety, P. A. (2007). Persecutory delusions and catastrophic worry in psychosis: developing the understanding of delusion distress and persistence. *Behaviour research and therapy*, 45(3), 523-537.
11. Freeman, D., Garety, P. A., Kuipers, E., Fowler, D., Bebbington, P. E., & Dunn, G. (2007). Acting on persecutory delusions: the importance of safety seeking. *Behaviour research and therapy*, 45(1), 89-99.
12. Moritz, S., Van Quaquebeke, N., & Lincoln, T. M. (2012). Jumping to conclusions is associated with paranoia but not general suspiciousness: a comparison of two versions of the probabilistic reasoning paradigm. *Schizophrenia research and treatment*, 2012.
13. Bronstein, M. V., Everaert, J., Castro, A., Joormann, J., & Cannon, T. D. (2019). Pathways to paranoia: Analytic thinking and belief flexibility. *Behaviour research and therapy*, 113, 18-24.
14. Freeman, D., Garety, P. A., Kuipers, E., Fowler, D., Bebbington, P. E., & Dunn, G. (2007). Acting on persecutory delusions: the importance of safety seeking. *Behaviour research and therapy*, 45(1), 89-99.
15. Murphy, P., Bentall, R. P., Freeman, D., O'Rourke, S., & Hutton, P. (2018). The paranoia as defense model of persecutory delusions: a systematic review and meta-analysis. *The Lancet Psychiatry*, 5(11), 913-929.

16. Valmaggia, L. R., Freeman, D., Green, C., Garety, P., Swapp, D., Antley, A., ... & Slater, M. (2007). Virtual reality and paranoid ideations in people with an 'at-risk mental state' for psychosis. *The British Journal of Psychiatry*, 191(S51), s63-s68.
17. McDonnell, J., Stahl, D., Day, F., McGuire, P., & Valmaggia, L. R. (2018). Interpersonal sensitivity in those at clinical high risk for psychosis mediates the association between childhood bullying victimization and paranoid ideation: a virtual reality study. *Schizophrenia research*, 192, 89-95.
18. Collip, D., Myin-Germeys, I., & Van Os, J. (2008). Does the concept of "sensitization" provide a plausible mechanism for the putative link between the environment and schizophrenia? *Schizophrenia bulletin*, 34(2), 220-225.
19. van Winkel, R., van Nierop, M., Myin-Germeys, I., & van Os, J. (2013). Childhood trauma as a cause of psychosis: linking genes, psychology, and biology. *The Canadian Journal of Psychiatry*, 58(1), 44-51.
20. Kapur, S., Mizrahi, R., & Li, M. (2005). From dopamine to salience to psychosis—linking biology, pharmacology and phenomenology of psychosis. *Schizophrenia research*, 79(1), 59-68.
21. Howes, O. D., Bose, S. K., Turkheimer, F., Valli, I., Egerton, A., Valmaggia, L. R., ... & McGuire, P. (2011). Dopamine synthesis capacity before onset of psychosis: a prospective [18F]-DOPA PET imaging study. *American Journal of Psychiatry*, 168(12), 1311-1317.
22. Howes, O. D., Kambeitz, J., Kim, E., Stahl, D., Slifstein, M., Abi-Dargham, A., & Kapur, S. (2012). The nature of dopamine dysfunction in schizophrenia and what this means for treatment: meta-analysis of imaging studies. *Archives of general psychiatry*, 69(8), 776-786.
23. Howes, O. D., McCutcheon, R., Owen, M. J., & Murray, R. M. (2017). The role of genes, stress, and dopamine in the development of schizophrenia. *Biological psychiatry*, 81(1), 9-20.
24. Schlier, B., Helbig-Lang, S., & Lincoln, T. M. (2016). Anxious but thoroughly informed? No jumping-to-conclusions bias in social anxiety disorder. *Cognitive Therapy and Research*, 40(1), 46-56.
25. Ellett, L., Allen-Crooks, R., Stevens, A., Wildschut, T., & Chadwick, P. (2013). A paradigm for the study of paranoia in the general population: The Prisoner's Dilemma Game. *Cognition and Emotion*, 27(1), 53–62.  
<https://doi.org/10.1080/02699931.2012.689757>
26. Savulich, G., Jeanes, H., Rossides, N., Kaur, S., Zacharia, A., Robbins, T. W., & Sahakian, B. J. (2018). Moral emotions and social economic games in paranoia. *Frontiers in psychiatry*, 9.
27. Haralanova, E., Haralanov, S., Beraldi, A., Möller, H. J., & Hennig-Fast, K. (2012). Subjective emotional over-arousal to neutral social scenes in paranoid schizophrenia. *European archives of psychiatry and clinical neuroscience*, 262(1), 59-68.
28. Williams, L. L. M., Das, P., Liddell, B. J., Olivieri, G., Peduto, A. S., David, A. S., ... & Harris, A. W. (2007). Fronto-limbic and autonomic disjunctions to negative emotion distinguish schizophrenia subtypes. *Psychiatry Research: Neuroimaging*, 155(1), 29-44.
29. Raihani, N. J., & Bell, V. (2017a). Conflict and cooperation in paranoia: a large-scale behavioral experiment. *Psychological Medicine*, pp. 1–11.  
<https://doi.org/10.1017/S0033291717003075>

30. Raihani, N. J., & Bell, V. (2017b). Paranoia and the social representation of others: A large-scale game theory approach. *Scientific Reports*, 7(1), 4544. <https://doi.org/10.1038/s41598-017-04805-3>
31. Greenburgh, A., Bell, V., & Raihani, N. (2018). Paranoia and conspiracy: group cohesion increases harmful intent attribution in the Trust Game. DOI: 10.31234/osf.io/mgzjr
32. Saalfeld, V., Ramadan, Z., Bell, V., & Raihani, N. J. (2018). Experimentally induced social threat increases paranoid thinking. *Royal Society Open Science*, 5(8), 180569.
33. Kahneman, D., Knetsch, J.L., Thaler, R.H. (1986) Fairness as a Constraint on Profit Seeking: Entitlements in the Market. *American Economic Review*. 76 (4): 728–741.
34. Berinsky, A. J., Huber, G. A., & Lenz, G. S. (2012). Evaluating online labor markets for experimental research: Amazon. com's Mechanical Turk. *Political analysis*, 20(3), 351-368.
35. Green, C. E. L., Freeman, D., Kuipers, E., Bebbington, P., Fowler, D., Dunn, G., & Garety, P. A. (2008). Measuring ideas of persecution and social reference: The Green et al. Paranoid Thought Scales (GPTS). *Psychological Medicine*, 38(1), 101–111.
36. Statham, V., Emerson, L. M., & Rowse, G. (2018). A Systematic Review of Self-Report Measures of Paranoia. *Psychological Assessment*. <https://doi.org/10.1037/pas0000645>
37. Grueber, C. E., Nakagawa, S., Laws, R. J., & Jamieson, I. G. (2011). Multimodel inference in ecology and evolution: Challenges and solutions. *Journal of Evolutionary Biology*. John Wiley & Sons, Ltd (10.1111). <https://doi.org/10.1111/j.1420-9101.2010.02210.x>
38. Burnham, K. P., & Anderson, D. R. (2004). Multimodel inference: Understanding AIC and BIC in model selection. *Sociological Methods and Research*. <https://doi.org/10.1177/0049124104268644>
39. Galipaud, M., Gillingham, M. A. F., David, M., & Dechaume-Moncharmont, F. X. (2014). Ecologists overestimate the importance of predictor variables in model averaging: A plea for cautious interpretations. *Methods in Ecology and Evolution*, 5(10), 983–991. <https://doi.org/10.1111/2041-210X.12251>
40. Barton, K. (2018). *Package “MuMIn” Title Multi-Model Inference*. Retrieved from <https://cran.r-project.org/web/packages/MuMIn/MuMIn.pdf>
41. Team, R. D. C., & R Development Core Team, R. (2016). R: A Language and Environment for Statistical Computing. *R Foundation for Statistical Computing*, 1(2.11.1), 409. <https://doi.org/10.1007/978-3-540-74686-7>
42. Wickham, H. (2016). *ggplot2: Elegant Graphics for Data Analysis*. New York: Springer. Retrieved from <https://cran.r-project.org/web/packages/ggplot2/citation.html>
43. Christensen, M. R. H. B. (2015). Package ‘ordinal’. *Stand*, 19, 2016.
44. Boyce, P., & Parker, G. (1989). Development of a scale to measure interpersonal sensitivity. *The Australian and New Zealand Journal of Psychiatry*, 23(3), 341–51.
45. Spielberger, C. D. (1989). *State-Trait Anxiety Inventory: Bibliography* (2nd ed.). Palo Alto, CA: Consulting Psychologists Press.
46. Meyer, T. J., Miller, M. L., Metzger, R. L., & Borkovec, T. D. (1990). Development and validation of the penn state worry questionnaire. *Behaviour research and therapy*, 28(6), 487-495.

47. Bates, D., Maechler, M., Bolker, B., Walker, S., Christensen, R. H. B., Singmann, H., ... & Grothendieck, G. (2011). Package 'lme4'. *Linear mixed-effects models using Eigen and S4 classes. R package version*, 1-1.
48. Goodrich, B., Gabry, J., Ali, I., & Brilleman, S. (2018) rstanarm: Bayesian applied regression modelling via Stan. *R Package Version 2.17.4*.
49. Mackowski, D., Ben-Shachar, M. S., Ludecke, D. (2019) Understand and Describe Bayesian Models and Posterior Distributions using bayestest. *RPackage*.
50. Kesting, M. L., Bredenkopf, M., Klenke, J., Westermann, S., & Lincoln, T. M. (2013). The impact of social stress on self-esteem and paranoid ideation. *Journal of behavior therapy and experimental psychiatry*, 44(1), 122-128.
51. Lincoln, T. M., Sundag, J., Schlier, B., & Karow, A. (2017). The relevance of emotion regulation in explaining why social exclusion triggers paranoia in individuals at clinical high risk of psychosis. *Schizophrenia bulletin*, 44(4), 757-767.
52. Sundag, J., Ascone, L., & Lincoln, T. M. (2018). The predictive value of early maladaptive schemas in paranoid responses to social stress. *Clinical psychology & psychotherapy*, 25(1), 65-75.
53. Kramer, I., Simons, C. J., Wigman, J. T., Collip, D., Jacobs, N., Derom, C., ... & Wichers, M. (2013). Time-lagged moment-to-moment interplay between negative affect and paranoia: new insights in the affective pathway to psychosis. *Schizophrenia bulletin*, 40(2), 278-286.
54. Barrantes-Vidal, N., Chun, C. A., Myin-Germeys, I., & Kwapil, T. R. (2013). Psychometric schizotypy predicts psychotic-like, paranoid, and negative symptoms in daily life. *Journal of Abnormal Psychology*, 122(4), 1077.
55. Schlier, B., Krkovic, K., Clamor, A., & Lincoln, T. M. (2019). Autonomic arousal during psychosis spectrum experiences: Results from a high-resolution ambulatory assessment study over the course of symptom on-and offset. *Schizophrenia Research*.
56. Veling, W., Pot-Kolder, R., Counotte, J., van Os, J., & van der Gaag, M. (2016). Environmental social stress, paranoia and psychosis liability: a virtual reality study. *Schizophrenia bulletin*, 42(6), 1363-1371.
57. Ellett, L., Freeman, D., & Garety, P. A. (2008). The psychological effect of an urban environment on individuals with persecutory delusions: the Camberwell walk study. *Schizophrenia research*, 99(1-3), 77-84.
58. Siegel, J. Z., Mathys, C., Rutledge, R. B., & Crockett, M. J. (2018). Beliefs about bad people are volatile. *Nature Human Behaviour*, 2(10), 750.
59. Freeman, D., Stahl, D., McManus, S., Meltzer, H., Brugha, T., Wiles, N., & Bebbington, P. (2012). Insomnia, worry, anxiety and depression as predictors of the occurrence and persistence of paranoid thinking. *Social psychiatry and psychiatric epidemiology*, 47(8), 1195-1203.
60. Freeman, D. (2016). Persecutory delusions: a cognitive perspective on understanding and treatment. *The Lancet Psychiatry*, 3(7), 685-692.
61. Sun, X., So, S. H. W., Chiu, C. D., Chan, R. C. K., & Leung, P. W. L. (2018). Paranoia and anxiety: A cluster analysis in a non-clinical sample and the relationship with worry processes. *Schizophrenia research*, 197, 144-149.
62. Startup, H., Pugh, K., Dunn, G., Cordwell, J., Mander, H., Černis, E., ... & Freeman, D. (2016). Worry processes in patients with persecutory delusions. *British Journal of Clinical Psychology*, 55(4), 387-400.

63. Freeman, D., Dunn, G., Startup, H., Pugh, K., Cordwell, J., Mander, H., ... & Kingdon, D. (2015). Effects of cognitive behaviour therapy for worry on persecutory delusions in patients with psychosis (WIT): a parallel, single-blind, randomised controlled trial with a mediation analysis. *The Lancet Psychiatry*, 2(4), 305-313.
64. Freeman, D, Slater, M, Bebbington, PE, Garety, PA, Kuipers, E, Fowler, D et al. (2003) Can virtual reality be used to investigate persecutory ideation? *The Journal of Nervous and Mental Disease* 191(8), 509–514.
65. Hartley, S., Haddock, G., e Sa, D. V., Emsley, R., & Barrowclough, C. (2014). An experience sampling study of worry and rumination in psychosis. *Psychological Medicine*, 44(8), 1605-1614.
66. Meisel, S. F., Garety, P. A., Stahl, D., & Valmaggia, L. R. (2018). Interpersonal processes in paranoia: a systematic review. *Psychological medicine*, 48(14), 2299-2312.
67. Freeman, D., Pugh, K., Vorontsova, N., Antley, A., & Slater, M. (2010). Testing the continuum of delusional beliefs: An experimental study using virtual reality. *Journal of abnormal psychology*, 119(1), 83.
68. Freeman, D, Gittins, M, Pugh, K, Antley, A, Slater, M and Dunn, G (2008) What makes one-person paranoid and another person anxious? The differential prediction of social anxiety and persecutory ideation in an experimental situation. *Psychological Medicine* 38(8), 1121–1132.
69. Freeman, D, Pugh, K, Antley, A, Slater, M, Bebbington, P, Gittins, M et al. (2008) Virtual reality study of paranoid thinking in the general population. *The British Journal of Psychiatry* 192(4), 258–263.
70. Green, CE, Freeman, D, Kuipers, E, Bebbington, P, Fowler, D, Dunn, G et al. (2011) Paranoid explanations of experience: a novel experimental study. *Behavioural and Cognitive Psychotherapy* 39(1), 21.
71. Peer, E., Brandimarte, L., Samat, S., & Acquisti, A. (2017). Beyond the Turk: Alternative platforms for crowdsourcing behavioral research. *Journal of Experimental Social Psychology*, 70, 153-163.
72. Crump, M. J., McDonnell, J. V., & Gureckis, T. M. (2013). Evaluating Amazon's Mechanical Turk as a tool for experimental behavioral research. *PloS one*, 8(3), e57410.
73. Hauser, D. J., & Schwarz, N. (2016). Attentive Turkers: MTurk participants perform better on online attention checks than do subject pool participants. *Behavior research methods*, 48(1), 400-407.