

# Detecting communicative intent in a computerised test of joint attention

Nathan Caruana<sup>Corresp., 1, 2, 3</sup>, Genevieve McArthur<sup>1, 2, 4</sup>, Alexandra Woolgar<sup>1, 2, 3</sup>, Jon Brock<sup>2, 4, 5</sup>

<sup>1</sup> Department of Cognitive Science, Macquarie University, Sydney, NSW, Australia

<sup>2</sup> ARC Centre of Excellence in Cognition and its Disorders, Sydney, NSW, Australia

<sup>3</sup> Perception in Action Research Centre, Sydney, NSW, Australia

<sup>4</sup> Centre for Atypical Neurodevelopment, Sydney, NSW, Australia

<sup>5</sup> Department of Psychology, Macquarie University, Sydney, NSW, Australia

Corresponding Author: Nathan Caruana

Email address: nathan.caruana@mq.edu.au

The successful navigation of social interactions depends on a range of cognitive faculties – including the ability to achieve joint attention with others to share information and experiences. We investigated the influence that intention monitoring processes have on gaze-following response times during joint attention. We employed a virtual reality task in which 16 healthy adults engaged in a collaborative game with a virtual partner to locate a target in a visual array. In the *Search* task, the virtual partner was programmed to engage in non-communicative gaze shifts in search of the target, establish eye contact, and then display a communicative gaze shift to guide the participant to the target. In the *NoSearch* task, the virtual partner simply established eye contact and then made a single communicative gaze shift towards the target (i.e., there were no non-communicative gaze shifts in search of the target). Thus, only the *Search* task required participants to monitor their partner's communicative intent before responding to joint attention bids. We found that gaze following was significantly slower in the *Search* task than the *NoSearch* task. However, the same effect on response times was not observed when participants completed non-social control versions of the *Search* and *NoSearch* tasks, in which the avatar's gaze was replaced by arrow cues. These data demonstrate that the intention monitoring processes involved in differentiating communicative and non-communicative gaze shifts during the *Search* task had a measurable influence on subsequent joint attention behaviour. The empirical and methodological implications of these findings for the fields of autism and social neuroscience will be discussed.

1

## Detecting Communicative Intent in a Computerised Test of Joint Attention

3

4

5 Nathan Caruana<sub>1 2 3</sub>, Genevieve McArthur<sub>1 2 4</sub>, Alexandra Woolgar<sub>1 2 3</sub>,  
6 & Jon Brock<sub>2 4 5</sub>

7

8

<sup>9</sup>ARC Centre of Excellence in Cognition and its Disorders, Australia

10

11

12

13

14

15

16

18

**ABSTRACT**

19 The successful navigation of social interactions depends on a range of cognitive faculties –  
20 including the ability to achieve joint attention with others to share information and experiences.  
21 We investigated the influence that intention monitoring processes have on gaze-following  
22 response times during joint attention. We employed a virtual reality task in which 16 healthy  
23 adults engaged in a collaborative game with a virtual partner to locate a target in a visual array.  
24 In the *Search* task, the virtual partner was programmed to engage in non-communicative gaze  
25 shifts in search of the target, establish eye contact, and then display a communicative gaze shift  
26 to guide the participant to the target. In the *NoSearch* task, the virtual partner simply established  
27 eye contact and then made a single communicative gaze shift towards the target (i.e., there were  
28 no non-communicative gaze shifts in search of the target). Thus, only the Search task required  
29 participants to monitor their partner's communicative intent before responding to joint attention  
30 bids. We found that gaze following was significantly slower in the Search task than the  
31 NoSearch task. However, the same effect on response times was not observed when participants  
32 completed non-social control versions of the Search and NoSearch tasks, in which the avatar's  
33 gaze was replaced by arrow cues. These data demonstrate that the intention monitoring processes  
34 involved in differentiating communicative and non-communicative gaze shifts during the Search  
35 task had a measurable influence on subsequent joint attention behaviour. The empirical and  
36 methodological implications of these findings for the fields of autism and social neuroscience  
37 will be discussed.

38       Joint attention is defined as the simultaneous coordination of attention between a social  
39 partner and an object or event of interest (Bruner, 1974; 1995). It is an intentional,  
40 communicative act. In the prototypical joint attention episode, one person *initiates* joint attention

41 (IJA) by pointing, turning their head, or shifting their eye gaze to intentionally guide their social  
42 partner to an object or event in the environment. The partner must recognise the intentional  
43 nature of this initiating behaviour and *respond* to that joint attention bid (RJA) by directing their  
44 attention to the cued location (Bruinsma, Koegel, & Koegel, 2004).

45 The ability to engage in joint attention is considered critical for the normal development of  
46 language and for navigating social interactions (Adamson, Bakeman, Deckner, & Romski, 2009;  
47 Charman, 2003; Dawson et al., 2004; Mundy, Sigman, & Kasari, 1990; Murray et al., 2008;  
48 Tomasello, 1995) and its developmental delay is a hallmark of autism spectrum disorders (Lord  
49 et al., 2000; Stone, Ousley, & Littleford, 1997). Yet despite its importance to both typical and  
50 atypical development, very little is known about the neurocognitive mechanisms of joint  
51 attention. By definition, joint attention involves an interaction between two individuals. The  
52 challenge for researchers, therefore, has been to develop paradigms that achieve the ecological  
53 validity of a dynamic, interactive, social experience, whilst at the same time maintaining  
54 experimental control.

55 In a recent functional magnetic resonance imaging (fMRI) study, Schilbach et al. (2010)  
56 investigated the neural correlates of joint attention using a novel virtual reality paradigm. During  
57 the scan, participants' eye-movements were recorded as they interacted with an anthropomorphic  
58 avatar. They were told that the avatar's gaze was controlled by a confederate outside the scanner  
59 also using an eye-tracking device. In fact, the avatar was controlled by a computer algorithm that  
60 responded to the participant's own eye-movements. On RJA trials (referred to as OTHER\_JA by  
61 Schilbach et al., 2010), the avatar looked towards one of three squares positioned around his face,  
62 and participants were instructed to respond by looking at the same square. Participants also  
63 completed IJA trials in which the roles were reversed.

64 Similar tasks have been used in other fMRI studies using either gaze-contingent avatars  
65 (Oberwelland et al., 2016) or live-video links to a real social partner (Redcay et al., 2012; Saito  
66 et al., 2010). Together, these interactive paradigms represent an important step towards an  
67 ecologically valid measure of joint attention. There is, however, a potentially important  
68 limitation of the tasks used in these studies: in each task, every trial involved a single  
69 unambiguously communicative eye-gaze cue. On RJA trials, the participant's partner would  
70 make a single eye-movement towards the target location and the participant knew they were  
71 required to respond to that isolated cue. This differs from real-life joint attention episodes, which  
72 are embedded within complex ongoing social interactions. In real life, responding to a joint  
73 attention bid requires that the individual first identifies the intentional nature of their partner's  
74 behaviour. That is, they must decide whether or not the cue is one they should follow. We refer  
75 to this component of joint attention as "intention monitoring".

76 In a recent fMRI study, we developed a novel joint attention task to better capture this  
77 intention monitoring process (Caruana, Brock & Woolgar, 2015). Following Schilbach et al.  
78 (2010), participants played a cooperative game with an avatar whom they believed to be  
79 controlled by a real person (referred to as "Alan"), but was actually controlled by a gaze-  
80 contingent algorithm. The participant and avatar were both allotted onscreen houses to search for  
81 a burglar (see Figure 1). On IJA trials, the participant found the burglar, made eye contact with  
82 the avatar, and then guided the avatar to the burglar by looking back at the house in which the  
83 burglar was hiding. On RJA trials, the participant found all of their allotted houses to be empty.  
84 Once Alan had finished searching his own houses, he would make eye contact with the  
85 participant before guiding them towards the house containing the burglar.

86  
87  
88  
89  
90  
91  
92  
93  
94



95 *Figure 1.* Experimental display showing the central avatar (“Alan”) and the six  
96 houses in which the burglar could be hiding. Gaze areas of interest (GAOIs),  
97 are represented by blue rectangles, and were not visible to participants.

98 The critical innovation of this task was the initial search phase. This provided a natural and  
99 intuitive context in which participants could determine, on each trial, their role as either the  
100 responder or initiator in a joint attention episode (previous studies had provided explicit  
101 instructions; e.g., Schilbach et al., 2010; Redcay et al., 2012). More importantly for current  
102 purposes, the RJA trials required participants to monitor their partner’s communicative  
103 intentions. During each trial, the avatar made multiple non-communicative eye-movements as he  
104 searched his own houses. The participant had to ignore these eye-movements and respond only  
105 to the communicative “initiating saccades” that followed the establishment of eye contact. This is  
106 consistent with genuine social interactions in which eye contact is used to signal one’s readiness  
107 and intention to communicate, particularly in the absence of verbal cues (Cary, 1978).

108 We compared the RJA trials in this new paradigm to non-social control trials (referred to as  
109 RJA) in which the eye gaze cue was replaced by a green arrow superimposed over the avatar's  
110 face. Analysis of saccadic reaction times revealed that participants were significantly slower (by  
111 approximately 209 ms) to respond to the avatar's eye gaze cue than they were to respond to the  
112 arrow (Caruana et al., 2015). We have since replicated this finding in an independent sample of  
113 adults and, intriguingly, found that the effect is exaggerated in a group of autistic individuals  
114 (Caruana, Stieglitz Ham et al., in press).

115 One explanation for these findings is that they reflect the intention monitoring aspects of  
116 RJA. Specifically, participants are slower to respond to eye gaze cues than arrows because it  
117 takes time to identify the cue as being an intentional and communicative bid to initiate joint  
118 attention. In the control condition, the arrow presents an unambiguous attention cue, and so the  
119 participant does not need to decide whether they should respond to it or not. The implication here  
120 is that intention monitoring is a cognitively demanding operation that requires time to complete  
121 and is manifest in the response times to eye gaze cues.

122 However, before reaching such a conclusion, it is important to consider a number of  
123 alternative explanations. For example, it may be that participants responded faster in the RJA  
124 condition because the large green arrow cue, which extended towards the target location,  
125 provided a more salient spatial cue than the avatar's eyes. It is also possible that the mere context  
126 of social interaction may influence the way participants approach the task. In particular, when  
127 individuals believe they are interacting with an intentional human agent, mirroring and  
128 mentalising mechanisms are automatically recruited which exert a top-down effect on the neural  
129 processes governing visual perception or attention (Wykowska, Wiese, Prosser, Müller &  
130 Hamed, 2014).

131        The aim of the current study, therefore, was to test the intention monitoring account more  
132   directly by manipulating the intention monitoring component of the RJA task whilst controlling  
133   for both the perceptual properties of the stimulus and the social nature of the task. To this end,  
134   we tested a new sample of participants using the same task but with one further manipulation. On  
135   half the trials, we eliminated the search phase of the task. Thus, on RJA trials, the avatar only  
136   made a single eye movement to the target to initiate joint attention, and participants knew  
137   unambiguously that they should follow it. The gaze cues in the ‘Search’ and ‘NoSearch’ versions  
138   of the task were identical and in both cases, participants believed they were interacting with  
139   another human. Thus, only the intention monitoring account predicts an effect of task (Search  
140   versus NoSearch) on response times. Participants also completed Search and NoSearch versions  
141   of the control (RJAc) condition. Because the arrow cue is unambiguous whether or not it is  
142   preceded by a search phase, we did not predict any difference in response times. In other words,  
143   a condition (Social vs. Control) by task (Search vs. NoSearch) interaction would indicate that  
144   response times to joint attention bids are influenced by the intention monitoring processes that  
145   precede true RJA behaviours.

## 146                      Method

### 147                      Ethical Statement

148        The study was approved by the Human Research Ethics Committee at Macquarie  
149   University (MQ; reference number: 5201200021). Participants received course credit for their  
150   time and provided written consent before participating.

### 151                      Participants

152        Sixteen right-handed adults with typical development, normal vision, and no history of  
153   neurological impairment participated in this study (3 female,  $M_{age} = 19.92$ ,  $SD = 1.03$ ).

### 154                      Stimuli

155 We employed an interactive paradigm that we had previously used to investigate the  
156 neural correlates of RJA and IJA (Caruana, et al., 2015). The stimuli comprised an  
157 anthropomorphic avatar face, generated using *FaceGen* (Singular Inversions, 2008), that  
158 subtended 6.5 degrees of visual angle. The avatar's gaze was manipulated so that it could be  
159 directed either at the participant or towards one of the six houses that were presented on the  
160 screen (see Figure 1). The houses were arranged in two horizontal rows above and below the  
161 avatar and each subtended four degrees of visual angle.

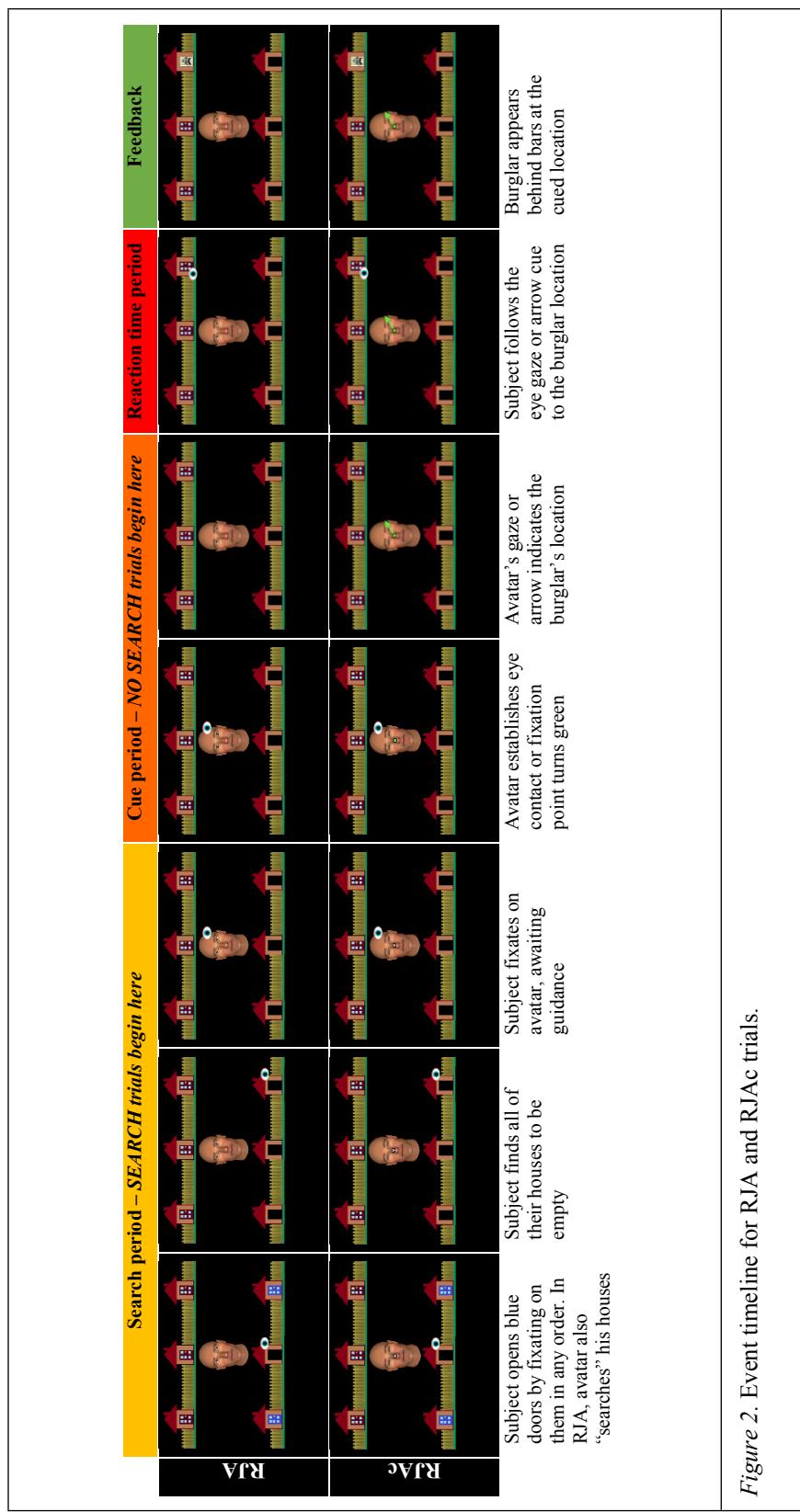


Figure 2. Event timeline for RJA and RJAc trials.

163           **Social Conditions (RJA and IJA).** Participants played a cooperative “Catch-the-  
164 Burglar” game with an avatar whom they believed was controlled by another person named  
165 “Alan” in a nearby eye tracking laboratory using live infrared eye tracking. In reality, a gaze-  
166 contingent algorithm controlled the avatar’s responsive behaviour (see Caruana et al., 2015 for a  
167 detailed description of this algorithm and a video demonstration of the task). The goal of the  
168 game was to catch a burglar that was hiding inside one of the six houses presented on the screen.  
169 Participants completed two versions of the social conditions (i.e., Search and NoSearch tasks)  
170 during separate blocks.

171           **Search Task.** This task was identical to the “Catch-the-Burglar” task employed in our  
172 previous work (e.g., Caruana et al., 2015). Each trial in the Search task began with a “search  
173 phase”. During this period, participants saw two rows of houses on a computer screen including  
174 a row of three blue doors and a row of three red doors. They were instructed to search the row of  
175 houses with blue doors while Alan searched the row of houses with red doors. Participants were  
176 told that they could not see the contents of Alan’s houses and that Alan could not see the  
177 contents of their houses. Whoever found the burglar first had to guide the other person to the  
178 correct location.

179           Participants searched the houses with blue doors in any order by fixating on them. Once a  
180 fixation was detected on a blue door, it opened to reveal either the burglar or an empty house. On  
181 some trials, only one or two blue doors were visible, whilst the remaining doors were already  
182 open. This introduced some variability in the order with which participants searched their houses  
183 that made Alan’s random search behaviour appear realistically unpredictable.

184           Once the participant fixated back on the avatar’s face, Alan was programmed to search 0-  
185 2 more houses and then make eye contact. This provided an interval in which participants could

186 observe Alan's non-communicative gaze behaviour as he completed his search. The onset  
187 latency of each eye movement made by Alan was jittered with a uniform distribution between  
188 500-1000 ms.

189 On RJA trials, participants discovered that all of their allotted houses were empty (Figure  
190 2, row 1), indicating that the burglar was hiding in one of Alan's houses. Once the participant  
191 fixated back on Alan's face, he searched 0-2 more houses in random order before establishing  
192 eye contact with the participant. Alan then initiated joint attention by directing his gaze towards  
193 one of his allotted houses. If the participant made a "responding saccade" and fixated the correct  
194 location, the burglar was captured.

195 On IJA trials, the participant found a burglar behind one of the blue doors. They were  
196 then required to fixate back on Alan's face, at which point the door would close again to conceal  
197 the burglar. Again, Alan was programmed to search 0-2 houses before looking straight at the  
198 participant. Once eye contact was established, participants could initiate joint attention by  
199 making an "initiating saccade" to fixate on the blue door that concealed the burglar. Alan was  
200 programmed to only respond to initiating saccades that followed the establishment of eye  
201 contact, and to follow the participant's gaze, irrespective of whether the participant fixated the  
202 correct house or not. Whilst performance on IJA trials was not of interest in the current study, the  
203 inclusion of this condition created a context for the collaborative search element of the task and  
204 allowed direct comparison with our previous studies in which participants alternated between  
205 initiating and responding roles.

206 When the participant made a responding or initiating saccade to the correct location, the  
207 burglar appeared behind prison bars to indicate that he had been successfully captured (e.g.,  
208 Figure 2, column 7). However, the burglar appeared in red at his true location, to indicate that he

209 had escaped, if participants (1) made a responding or initiating saccade to an incorrect location,  
210 (2) took longer than three seconds to make a responding or initiating saccade, or (3) spent more  
211 than three seconds looking away from task-relevant stimuli (i.e., Alan and houses). Furthermore,  
212 trials were terminated if the participants took longer than three seconds to begin searching their  
213 houses at the beginning of the trial. On these trials, red text reading “Failed Search” appeared on  
214 the screen to provide feedback.

215         **NoSearch Task.** This version of the task was identical to the Search task except that the  
216 search phase in each trial was removed. In IJA trials, all but one house was visibly empty (i.e.,  
217 the door was open and no burglar was present), and participants were instructed that if they saw a  
218 blue door in their allotted row of houses that the burglar would be “hiding” behind it. In RJA  
219 trials, all of the houses were visibly empty. For both IJA and RJA trials, Alan’s eyes would be  
220 closed at the beginning of the trial, and then open after 500-1000 ms (jittered with a uniform  
221 distribution) so that he was looking at the participant. Alan would then wait to be guided on IJA  
222 trials. On RJA trials, Alan shifted his gaze to guide the participant after a further 500-1000 ms,  
223 provided that eye contact had been maintained. Thus, in both the Search and NoSearch tasks,  
224 Alan made eye contact with the participant before guiding them to the burglar on RJA trials.  
225 Therefore, the perceptual properties of the gaze cue itself were identical between tasks, but the  
226 NoSearch task removed the requirement to use the eye contact cue to identify communicative  
227 gaze shifts.

228         **Control Conditions (RJAc and IJAc).** For each of the social conditions in both versions  
229 of the task, we employed a control condition that was closely matched on non-social task  
230 demands (e.g., attentional orienting, oculomotor control). In these conditions (RJAc and IJAc),  
231 participants were told that they would play the game without Alan, whose eyes remained closed

232 during the trial. Participants were told that the stimuli presented on the computer screen in these  
233 trials were controlled by a computer algorithm. In the Search task, a grey fixation point was  
234 presented over the avatar's nose until the participant completed their search and fixated upon it.  
235 After a short delay, the fixation point turned green (analogous to the avatar making eye contact).  
236 From this point onwards, the Search and NoSearch tasks were identical. On IJAc trials, the green  
237 fixation point was the cue to saccade towards the burglar location. On RJA trials a green arrow  
238 subtending three degrees of visual angle cued the burglar's location (analogous to Alan's guiding  
239 gaze; see Caruana et al., 2015 for a video with example trials from each condition).

#### 240 Procedure

241 **Joint attention task.** The experiment was presented using *Experiment Builder* 1.10.165  
242 (SR Research, 2004). Participants completed four blocks, each comprising 108 trials: two blocks  
243 involved the Search task, and another two blocks involved the NoSearch task. Search and  
244 NoSearch block pairs were presented consecutively, however their order was counterbalanced  
245 across participants. Within each pair of Search and NoSearch blocks, one block required the  
246 participant to monitor the upper row of houses, and the other required them to monitor the lower  
247 row of houses.

248 Each block comprised 27 trials from each condition (i.e., RJA, RJA, IJA, IJA). Social  
249 (RJA, IJA) and control (RJA, IJA) trials were presented in clusters of six trials throughout  
250 each block. Each cluster began with a cue lasting 1000 ms that was presented over the avatar  
251 stimulus and read "Together" for a social cluster and "Alone" for a control cluster. Trial order  
252 randomisation was constrained to ensure that the location of the burglar, the location of blue  
253 doors, and the number of gaze shifts made by the avatar were matched within each block and  
254 condition.

255 After playing the interactive game, and consistent with our previous studies employing  
256 this paradigm, a post-experimental interview was conducted in which participants were asked to  
257 rate their subjective experience during the task (cf. Caruana et al., 2015; Caruana, Ham et al., in  
258 pres). Full details on the assessment of subjective experiences and relevant findings are provided  
259 in Supplementary Material 2.

260 **Eye tracking.** Eye-movements from the right eye only were recorded with a sampling  
261 rate of 500Hz using a desktop-mounted EyeLink 1000 Remote Eye-Tracking System (SR  
262 Research Ltd., Ontario, Canada). Head movements were stabilised using a chinrest. We  
263 conducted an eye tracking calibration using a 9-point sequence at the beginning of each block.  
264 Seven gaze areas of interest (GAOIs) over the houses and avatar stimulus were used by our gaze-  
265 contingent algorithm (see Caruana et al., 2015 for details). A recalibration was conducted if the  
266 participant made consecutive fixations on the borders or outside the GAOIs. Trials requiring a  
267 recalibration were excluded from all analyses.

## 268 Scores

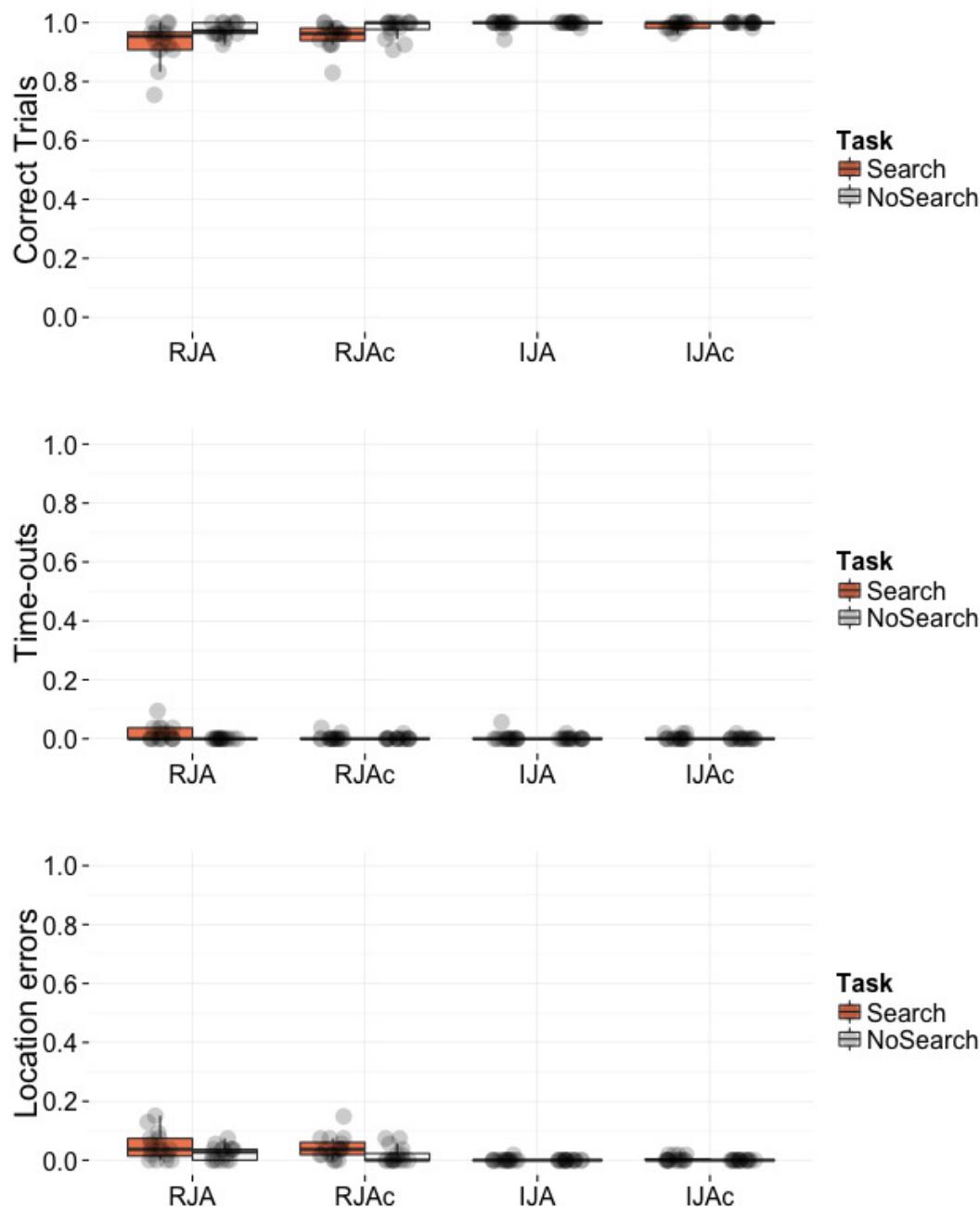
269 **Accuracy.** We calculated the proportion of trials where the participant succeeded in  
270 catching the burglar in each condition (i.e., RJA and RJAc) for each task separately (i.e., Search  
271 and NoSearch). We excluded from the accuracy analysis any trials that required a recalibration  
272 or (in the Search task) were failed due to an error during the search phase.

273 **Saccadic Reaction Times.** For correct trials, we measured the latency (in ms) between  
274 the presentation of the gaze cue (for RJA trials) or the arrow cue (for RJAc trials), and the onset  
275 of the participant's responding saccade towards the burglar location (see Figure 2, Reaction time  
276 period).

## 277 Statistical Analyses

278 Saccadic reaction times were analysed via repeated measures Analysis of Variance  
279 (ANOVA) using the ezANOVA (ez) package in R (Lawrence, 2013), reporting the generalised  
280 eta squared ( $\eta^2_G$ ) measure of effect size. Significant task\*condition interactions effects were  
281 followed-up with Welch's two sample unequal variances t-tests (Welch, 1947). As in our  
282 previous studies, we report analyses of the mean reaction time, having excluded trials with  
283 reaction times less than 150 ms as these are typically considered to be anticipatory responses.  
284 Trials timed out after 3000 ms providing a natural upper limit to reaction times. Full syntax and  
285 output for this analysis can be found in the Rmarkdown document (Supplementary Material 1).  
286 The RMarkdown also provides complementary ANOVAs of the mean and median of the  
287 untrimmed data, as well as a mixed random effects analysis using the lme4 R package (Bates,  
288 2005). The results of all analyses are consistent in terms of the predicted interaction between task  
289 and condition. A significance criterion of  $p < 0.05$  was used for all analyses.

291



292

293 *Figure 3.* Box plots displaying the proportion of correct trials, proportion of time-out errors, and proportion of  
294 location errors. Data points represent individual participant means.

295

296

## Results

297 As depicted in Figure 3, participants performed at close-to-ceiling levels in terms of the  
298 trials successfully completed across all conditions. Of the small number of errors made, the  
299 majority were Location Errors in the RJA and RJAc conditions, whereby participants looked first  
300 to an incorrect location (house) rather than following the avatar's gaze to the burglar location.  
301 Given the low number of errors, we do not report statistical analyses of accuracy or errors.

302 Figure 4 shows mean saccadic reaction times for correct trials in the RJA and RJAc  
303 conditions. Participants were significantly slower on the Search task than the NoSearch task  
304 (main effect of task ( $F(1,15) = 11.07, p = .005, \eta^2_G = 0.13$ ). They were also significantly slower to  
305 respond on RJA trials than RJAc trials overall (main effect of condition,  $F(1,15) = 98.75, p <$   
306  $.001, \eta^2_G = 0.57$ ). Importantly, there was also a significant task\*condition interaction ( $F(1,15) =$   
307  $43.86, p < .001, \eta^2_G = 0.18$ ), indicating a larger effect of task in the RJA condition. This  
308 interaction was also present in all re-analyses of the data (see Supplementary Material 1).

309 Follow-up paired t-tests revealed that responses to social gaze were significantly  
310 slower in the Search task than the NoSearch task ( $t(15) = 4.82, p < .001$ ), whereas response times  
311 to arrow cues did not significantly differ between the two versions of the control task ( $t(15) = -$   
312  $0.85, p = .411$ ). Consistent with our previous studies, there was an effect of condition (i.e.,  
313 slower responses to gaze cues than arrow cues) for the Search ( $t(15) = 9.31, p < .001$ ) task.  
314 However, there was also a significant (albeit smaller) effect of condition for the NoSearch task  
315 ( $t(15) = 8.51, p < .001$ ).

316

317

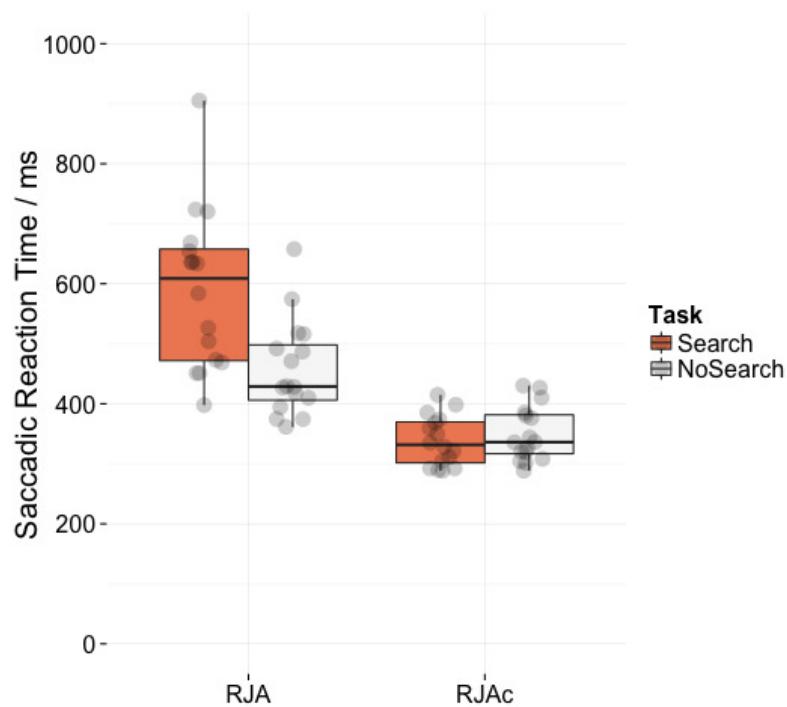


Figure 4. Box plots displaying saccadic reaction times in RJA and RJAc conditions, separated by task (i.e., Search, NoSearch). Data points represent individual participant means.

318

**Discussion**

319 One of the main challenges facing social neuroscience – and the investigation of joint  
320 attention in particular – is the need to achieve ecological validity whilst maintaining  
321 experimental control. During genuine joint attention experiences, our social cognitive faculties  
322 are engaged whilst we are immersed in complex interactions consisting of multiple social cues  
323 with the potential for communication. A critical but neglected aspect of joint attention is the  
324 requirement to identify those cues that are intended to be communicative. In the specific case of  
325 eye gaze cues, the responder must differentiate gaze shifts that signal an intentional joint  
326 attention bid from other, non-communicative gaze shifts. The results of the current study indicate  
327 that this intention monitoring process has a measurable effect on responsive joint attention  
328 behaviour.

329 The Search version of our Catch-the Burglar task was identical to that used in our  
330 previous studies. In the social (RJA) condition, participants found all of their houses to be empty,  
331 waited for their partner, Alan, to complete his search, make eye contact, and then guide them to  
332 the burglar's location. We replicated our previous finding (Caruana et al., 2015; Caruana et al., in  
333 press) that participants were slower to respond in this condition than in the matched control  
334 (RJAc) condition in which the avatar's eye gaze cue was replaced by an arrow.

335 The critical innovation of the current study was the addition of a NoSearch condition in  
336 which the same gaze and arrow cues were used but the joint attention episode was not preceded  
337 by a search phase, thereby removing the intention monitoring component of the RJA condition.  
338 Participants were still slower to respond to eye gaze than to arrow cues, suggesting that the  
339 previously identified difference between RJA and RJAc in the Search condition is not entirely  
340 attributable to intention monitoring. As discussed earlier, it is possible that differences in the

341 perceptual salience of the arrow cue might help contribute this effect. Alternatively, participants  
342 may be affected by the presence of a social partner. The current study was designed to control for  
343 such factors and it is not possible to determine which if either of these explanations is correct.

344 The important finding was the task by condition interaction. This arose because the  
345 magnitude of the condition effect was significantly reduced in the NoSearch version of the task.  
346 This cannot be explained in terms of perceptual salience or social context, because these were  
347 identical across Search and NoSearch tasks.

348 The interaction can also be viewed by contrasting the effect of task (Search vs.  
349 NoSearch) for the two different conditions (RJA or RJA). In the RJA condition, participants  
350 were significantly faster to respond to the eye gaze cue when the search phase was removed. The  
351 search phase required the participant and their virtual partner to make multiple non-  
352 communicative eye-movements prior to the joint attention episode. Participants therefore had to  
353 differentiate between eye-movements made by the avatar that signalled a communicative joint  
354 attention bid and those that were merely a continuation of their search. In the NoSearch task,  
355 every eye-movement made by the avatar was communicative, thereby removing the requirement  
356 to monitor his communicative intent, enabling faster response times.

357 Importantly, there was no effect of task (Search vs. NoSearch) for the RJA condition.  
358 This allows us to discount a number of alternative explanations for the task effect in the RJA  
359 condition. For example, it could be argued that the slower responses in the Search task reflected  
360 differences in the timing of the stimulus presentation (e.g., the delay between participants  
361 fixating on the avatar and the avatar making his guiding saccade). However, the timing of the  
362 stimuli were programmed to be identical in the corresponding RJA and RJA conditions of the  
363 Search and NoSearch tasks, so any effect of stimulus timing should have been evident in both

364 conditions. Another plausible explanation is that participants were slower in the Search task  
365 because this required them to switch from searching for the burglar to responding to the avatar  
366 on each trial. But again this applied equally to the RJA and RJAc conditions, so it cannot explain  
367 the task by condition interaction.

368 In short, the observed interaction between task and condition is entirely consistent with  
369 our intention monitoring account and cannot be explained in terms of the perceptual salience of  
370 different cues, the task's social context, the timing of the stimulus presentation, or the  
371 requirement to switch between searching and responding.

372 The current data provide insights into our other recent findings in studies using the  
373 Search version of our interactive task. In one study, we used fMRI to investigate the neural  
374 correlates of RJA (Caruana et al., 2015). By contrasting activation in the RJA (eye gaze) and  
375 RJAc (arrow) conditions, we identified a broad frontotemporoparietal network including the  
376 right temporo-parietal junction and right inferior frontal lobe. These brain regions are strongly  
377 associated with aspects of social cognition including mentalising (e.g., Saxe & Kanwisher, 2003)  
378 and predicting another's actions (Danckert et al., 2002; Hamilton & Grafton, 2008) but have not  
379 been previously linked to RJA (cf. Redcay et al., 2012; Schilbach et al., 2010). The tasks used in  
380 previous studies of RJA were similar to the current NoSearch task. As such, they would not have  
381 captured the intention monitoring processes involved in RJA, perhaps explaining the  
382 discrepancy with our fMRI study. Future neuroimaging studies of joint attention could employ  
383 the current study's design, and compare activation observed during the Search and NoSearch  
384 task. If our interpretation is correct then removing the search component should reduce the  
385 involvement of temporoparietal and inferior frontal regions in the RJA condition.

386 In another study (Caruana, Stieglitz Ham, et al., in press), we investigated joint attention  
387 in adults with autism. Observational studies of real-life interactions provide overwhelming  
388 evidence that joint attention impairments are a core feature of autism (Charman et al., 1997;  
389 Dawson et al., 2004; Loveland & Landry, 1986; Mundy et al., 1990; Osterling, Dawson, &  
390 Munson, 2002; Wong & Kasari, 2012). However, previous computer-based experimental studies  
391 of joint attention have largely failed to find consistent evidence of gaze following difficulties  
392 (Leekam, 2015; Nation & Penny, 2008). One possible explanation for this is that autistic  
393 individuals have an underlying difficulty in understanding the social significance or  
394 communicative intentions conveyed by eye contact (cf. Böckler, Timmermans, Sebanz, Vogeley,  
395 & Schilbach, 2014; Senju & Johnson, 2009) that is not captured by the tasks used in previous  
396 studies of autism. In contrast to these studies, we did find evidence of impairment: autistic adults  
397 made more errors and were slower to respond than control participants in the RJA condition  
398 despite showing no impairment on the control condition. Future studies involving individuals  
399 with autism and the Search and NoSearch versions of our task would clarify this issue further.

400 Another avenue for further research using this task would be to investigate the  
401 development of joint attention in young children. Infants begin responding to and initiating joint  
402 attention bids within the first year of life (Mundy et al., 2007), but virtual reality tasks provide  
403 the sensitivity to investigate the developmental changes in the speed and efficiency of joint  
404 attention engagement in later development. Finally, it would be of interest to investigate sex  
405 differences in performance. Studies of infants (Saxon & Reilly, 1999) and school-aged children  
406 (Gavrilov, Rotem, Ofek & Geva, 2012) have found that females exhibit increased joint attention  
407 behaviours compared to their male peers, although it is unclear to what extent these differences  
408 reflect underlying differences in competence as opposed to motivation. A limitation of the

409 current study is that only three participants were female. However, future studies with larger  
410 samples would allow systematic investigation of sex differences in joint attention performance at  
411 multiple points across development.

412 **Summary**

413 In everyday joint attention episodes, a critical aspect of responding to joint attention bids  
414 is the ability to discern which social cues have communicative intent and which do not. The  
415 results of the current study indicate that this intention monitoring component has a measureable  
416 effect on responding behaviour. Moreover, this component can be isolated by contrasting joint  
417 attention episodes occurring in the context of a realistically complex social interaction versus a  
418 simplified context in which each cue is unambiguously communicative. The clear differences in  
419 performance on the Search and NoSearch versions of our task highlight the importance of  
420 striving for ecological validity in studies of social cognition (cf. Schilbach et al., 2013). The  
421 results also demonstrate the potential of our task for investigating the different components of  
422 joint attention in typically developing children and in clinical populations associated with  
423 atypical social cognition.

424

**References**

- 425 Adamson, L. B., Bakeman, R., Deckner, D. F., & Romski, M. A. (2009). Joint engagement and  
426 the emergence of language in children with autism and Down syndrome. *Journal of*  
427 *Autism and Developmental Disorders*, 39(1), 84-96. doi: 10.1007/s10803-008-0601-7
- 428 Bayliss, A. P., di Pellegrino, G., & Tipper, S. P. (2005). Sex differences in eye gaze and  
429 symbolic cueing of attention. *Quarterly Journal of Experimental Psychology*, 58, 1-20.
- 430 Bates DM. 2005. Fitting linear mixed models in R. *R News* 5:27-30
- 431 Böckler, A., Timmermans, B., Sebanz, N., Vogeley, K., & Schilbach, L. (2014). Effects of  
432 observing eye contact on gaze following in high-functioning autism. *Journal of Autism and*  
433 *Developmental Disorders*, 44(7), 1651-1658. doi: 10.1007/s10803-014-2038-5
- 434 Bruinsma, Y., Koegel, R. L., & Koegel, L. K. (2004). Joint attention and children with autism: A  
435 review of the literature. *Mental Retardation & Developmental Disabilities Research*  
436 *Reviews*, 10(3), 169-175. doi: 10.1002/mrdd.20036
- 437 Bruner, J. S. (1974). From communication to language, A psychological perspective. *Cognition*,  
438 3(3), 255-287. doi: 10.1016/0010-0277(74)90012-2
- 439 Bruner, J. S. (1995). From joint attention to the meeting of minds: an introduction. In C. Moore  
440 & P. J. Dunham (Eds.), *Joint attention: its origins and role in development*. (pp. 1-14).  
441 Hillsdale, NJ: Lawrence Erlbaum Associates.
- 442 Caruana, N., Brock, J., & Woolgar, A. (2015). A frontotemporoparietal network common to  
443 initiating and responding to joint attention bids. *NeuroImage*, 108, 34-46. doi:  
444 10.1016/j.neuroimage.2014.12.041

- 445 Caruana, N., de Lissa, P., & McArthur, G. (2016). Beliefs about human agency influence the  
446 neural processing of gaze during joint attention. *Social Neuroscience*,  
447 doi:10.1080/17470919.2016.1160953
- 448 Caruana, N., Stieglitz Ham., H., Brock, J., Woolgar, A., Palermo, R., Kloth, N., & McArthur, G.  
449 (in press). Joint Attention Difficulties in Adults with Autism. *Autism*.
- 450 Cary, M. S. (1978). The role of gaze in the initiation of conversation. *Social Psychology*, 41(3),  
451 269-271. doi: 10.2307/3033565
- 452 Charman, T. (2003). Why is joint attention a pivotal skill in autism? *Philosophical Transactions  
453 Royal Society London Biological Sciences*, 358, 315-324. doi: 10.1098/rstb.2002.1199
- 454 Charman, T., Swettenham, J., Baron-Cohen, S., Cox, A., Baird, G., & Drew, A. (1997). Infants  
455 with autism: An investigation of empathy, pretend play, joint attention, and imitation.  
456 *Developmental Psychology*, 33(5), 781-789. doi: 10.1037/0012-1649.33.5.781
- 457 Danckert J, Ferber S, Doherty T, Steinmetz H, Nicolle D, Goodale MA (2002) Selective, non-  
458 lateralized impairment of motor imagery following right parietal damage. *Neurocase* 8,  
459 194–204
- 460 Dawson, G., Toth, K., Abbott, R., Osterling, J., Munson, J. A., Estes, A., & Liaw, J. (2004).  
461 Early social attention impairments in autism: Social orienting, joint attention, and  
462 attention to distress. *Developmental Psychology*, 40(2), 271-283. doi:  
463 <http://dx.doi.org/10.1037/0012-1649.40.2.271>
- 464 Gavrilov, Y., Rotem, S., Ofek, R., & Geva, R. (2012). Socio-cultural effects on children's  
465 initiation of joint attention. *Frontiers in Human Neuroscience*, 6, 286.
- 466 Hamilton AFC, de, Grafton ST. 2008. Action outcomes are represented in human inferior  
467 frontoparietal cortex. *Cerebral Cortex*, 18 (5), 1160–1168.

- 468 Lawrence MA (2013). ez: Easy analysis and visualization of factorial experiments. R package  
469 version 4.2-2, URL <http://CRAN.R-project.org/package=ez>.
- 470 Leekam, S. (2015). Social cognitive impairment and autism: what are we trying to explain?  
471 *Philosophical Transactions of the Royal Society of London B: Biological Sciences*, 371,  
472 1686.
- 473 Lord, C., Risi, S., Lambrecht, L., Cook, E.H., Leventhal, B.L., DiLavore, P.C., Pickles, A.,  
474 Rutter, M., 2000. The autism diagnostic observation schedule—generic: a standard  
475 measure of social and communication deficits associated with the spectrum of autism.  
476 *Journal of Autism and Developmental Disorders*, 30(3), 205–223. doi:  
477 10.1023/A:1005592401947
- 478 Loveland, K. A., & Landry, S. H. (1986). Joint attention and language in autism and  
479 developmental language delay. *Journal of Autism and Developmental Disorders*, 16(3),  
480 335-349. doi: 10.1007/BF01531663
- 481 Mundy, P., Sigman, M., & Kasari, C. (1990). A longitudinal study of joint attention and  
482 language development in autistic children. *Journal of Autism and Developmental  
483 Disorders*, 20(1), 115-128. doi: 10.1007/bf02206861
- 484 Mundy, P., Block, J., Vaughan Van Hecke, A., Delgado, C., Parlade, M., Pomeras, Y. (2007).  
485 Individual differences in the development of joint attention in infancy. *Child  
486 Development*, 78, 938–954.
- 487 Murray, D. S., Creaghead, N. A., Manning-Courtney, P., Shear, P. K., Bean, J., & Prendeville, J.  
488 A. (2008). The relationship between joint attention and language in children with autism  
489 spectrum disorders. *Focus on Autism & Other Developmental Disabilities*, 23(1), 5-14.  
490 doi: 10.1177/1088357607311443

- 491 Nation, K., & Penny, S. (2008). Sensitivity to eye gaze in autism: Is it normal? Is  
492 it social? *Development and Psychopathology*, 20(01), 79-97. doi:  
493 doi:10.1017/S0954579408000047
- 494 Oberwelland, E., Schilbach, L., Barisic, I., Krall, S. C., Vogeley, K., Fink, G. R., Herpertz-  
495 Dahlmann, B., Konrad, K., Schulte-Rüther, M. (2016). Look into my eyes: Investigating  
496 joint attention using interactive eye-tracking and fMRI in a developmental sample.  
497 *NeuroImage*, in press.
- 498 Osterling, J., Dawson, G., & Munson, J. A. (2002). Early recognition of 1-year-old infants with  
499 autism spectrum disorder versus mental retardation. *Development and Psychopathology*,  
500 14(02), 239-251. doi: doi:10.1017/S0954579402002031
- 501 Redcay, E., Dodell-Feder, D., Mavros, P. L., Kleiner, A. M., Pearrow, M. J., Triantafyllou, C.,  
502 Gabrieli, J. D., & Saxe, R. (2012). Atypical brain activation patterns during a face-to-face  
503 joint attention game in adults with autism spectrum disorder. *Human Brain Mapping*, 34,  
504 2511-2523. doi: 10.1002/hbm.22086.
- 505 Saito, D. N., Tanabe, H. C., Izuma, K., Hayashi, M. J., Morito, Y., Komeda, H., . . . Sadato, N.  
506 (2010). "Stay Tuned": inter-individual neural synchronization during mutual gaze and  
507 joint attention. *Frontiers in Integrative Neuroscience*, 4, 127. doi:  
508 10.3389/fnint.2010.00127
- 509 Saxe, R., & Kanwisher, N. (2003). People thinking about thinking people: The role of the  
510 temporo-parietal junction in "theory of mind". *NeuroImage*, 19(4), 1835–1842. doi:  
511 [http://dx.doi.org/10.1016/S1053-8119\(03\)00230-1](http://dx.doi.org/10.1016/S1053-8119(03)00230-1).
- 512 Saxon, Terrill F., and John T. Reilly. (1999). Joint Attention and Toddler Characteristics: Race,  
513 Sex and Socioeconomic Status. *Early Child Development and Care*, 149(1), 59-69.

- 514 Schilbach, L., Wilms, M., Eickhoff, S. B., Romanzetti, S., Tepest, R., Bente, G., Shah N. J.,  
515 Fink, G. R., & Vogeley, K. (2010). Minds made for sharing: Initiating joint attention  
516 recruits reward-related neurocircuitry. *Journal of Cognitive Neuroscience*, 22(12), 2702-  
517 2715. doi: 10.1162/jocn.2009.21401.
- 518 Schilbach, L., Timmermans, B., Reddy, V., Costall, A., Bente, G., Schlicht, T., & Vogeley, K.  
519 (2013). Toward a second-person neuroscience. *Behavioral and Brain Sciences*, 36(4),  
520 393-414. doi: 10.1017/S0140525X12000660.
- 521 Senju, A., & Johnson, M. H. (2009). The eye contact effect: mechanisms and development.  
522 *Trends in Cognitive Sciences*, 13(3), 127-134. doi:  
523 <http://dx.doi.org/10.1016/j.tics.2008.11.009>
- 524 Singular Inversions. (2008). FaceGen Modeller (Version 3.3) [Computer Software]. Toronto,  
525 ON: Singular Inversions.
- 526 SR Research. (2004). Experiment Builder (Version 1.10.165). Ontario.
- 527 Stone, W. L., Ousley, O. Y., & Littleford, C. D. (1997). Motor imitation in young children with  
528 autism: what's the object? *Journal of Abnormal Child Psychology*, 25, 475-485. doi:  
529 10.1023/A:1022685731726
- 530 Tomasello, M. (1995). Joint attention as social cognition. In C. Moore & P. J. Dunham (Eds.),  
531 *Joint Attention: Its Origins and Role in Development*. Hillsdale: Lawrence Erlbaum  
532 Associates.
- 533 Welch, B. L. (1947). The generalization of 'student's' problem when several different population  
534 variances are involved. *Biometrika*, 34, 28-35.

- 535 Wong, C., & Kasari, C. (2012). Play and joint attention of children with autism in the preschool  
536 special education classroom. *Journal of Autism and Developmental Disorders*, 42(10),  
537 2152-2161. doi: 10.1007/s10803-012-1467-2
- 538 Wykowska, A., Wiese, E., Prosser, A., Müller, H. J., & Hamed, S. B. (2014). Beliefs about the  
539 minds of others influence how we process sensory information. *PLoS ONE*, 9(4), e94339.  
540 doi:10.1371/journal.pone.0094339