

# 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

2

**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

<sub>1</sub> Department of Cognitive Science, MACQUARIE UNIVERSITY, Sydney, Australia

9

<sub>2</sub>ARC Centre of Excellence in Cognition and its Disorders, Australia

10

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

11

<sub>4</sub>Centre for Atypical Neurodevelopment, MACQUARIE UNIVERSITY, Sydney, Australia

12

<sub>5</sub> Department of Psychology, MACQUARIE UNIVERSITY, Sydney, Australia

13

**Corresponding Author:**

15 Dr. Nathan Caruana

16 Email: [nathan.caruana@mq.edu.au](mailto:nathan.caruana@mq.edu.au)

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, & Ronski, 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

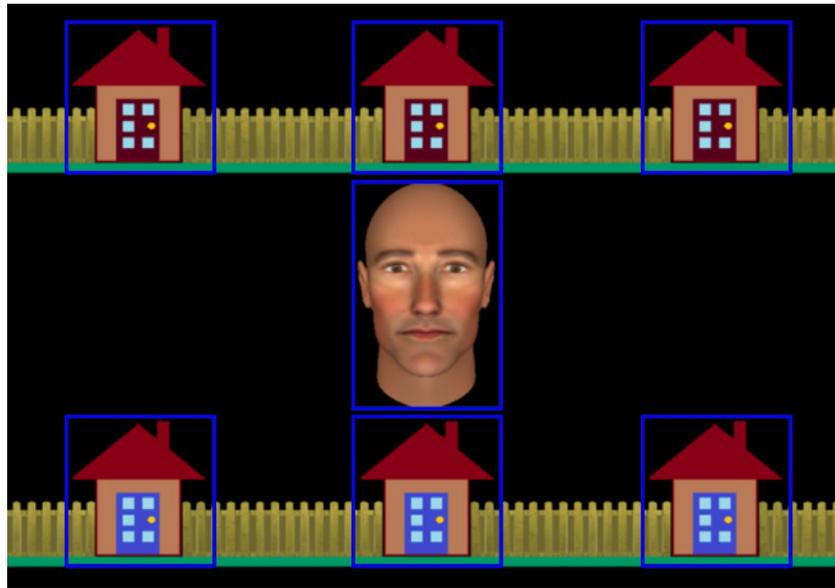
90

91

92

93

94



95

96

97

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

98

99

100

101

102

103

104

105

106

107

The critical innovation of this task was the initial search phase. This provided a natural and intuitive context in which participants could determine, on each trial, their role as either the responder or initiator in a joint attention episode (previous studies had provided explicit instructions; e.g., Schilbach et al., 2010; Redcay et al., 2012). More importantly for current purposes, the RJA trials required participants to monitor their partner’s communicative intentions. During each trial, the avatar made multiple non-communicative eye-movements as he searched his own houses. The participant had to ignore these eye-movements and respond only to the communicative “initiating saccades” that followed the establishment of eye contact. This is consistent with genuine social interactions in which eye contact is used to signal one’s readiness 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 RJAc) in which the eye gaze cue was replaced by a green arrow superimposed over the avatar's  
110 face. Analysis of saccadic reactions 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). This effect was surprising – previous studies have shown that gaze  
113 cues often engender rapid and reflexive attention shifts (see Frischen, Bayliss, and Tipper, 2007),  
114 but that would predict faster rather than slower responses to gaze cues. Nevertheless, we have  
115 since replicated this finding in an independent sample of adults and, intriguingly, found that the  
116 effect is exaggerated in a group of autistic individuals (Caruana, Stieglitz Ham et al., in press).

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

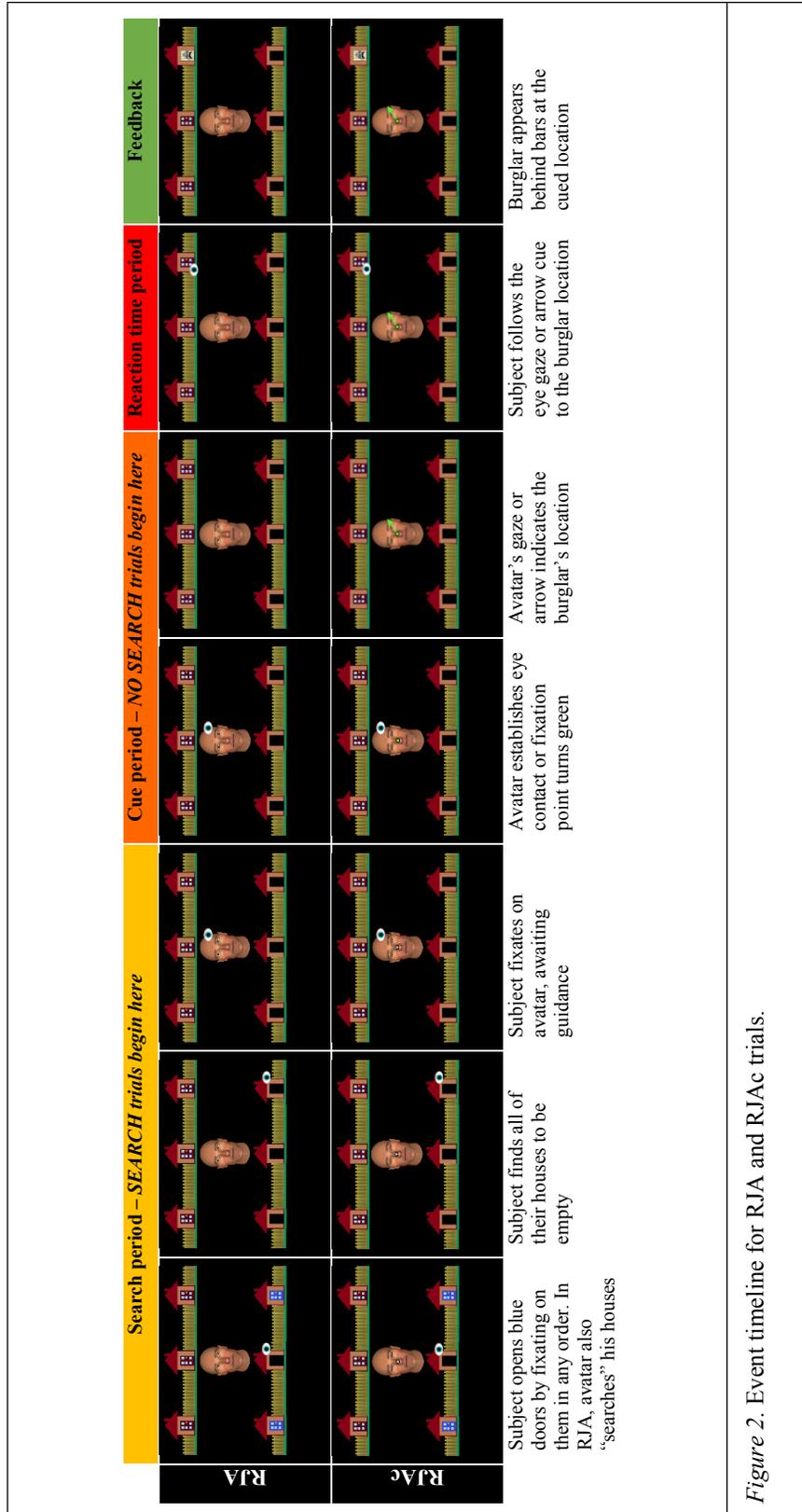
124 However, before reaching such a conclusion, it is important to consider a number of  
125 alternative explanations. For example, it may be that participants responded faster in the RJAc  
126 condition because the large green arrow cue, which extended towards the target location,  
127 provided a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



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

### 156 **Stimuli**

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



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

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

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

186           Once the participant fixated back on the avatar’s face, Alan was programmed to search 0-  
187 2 more houses and then make eye contact. This provided an interval in which participants could

188 observe Alan's non-communicative gaze behaviour as he completed his search. The onset  
189 latency of each eye movement made by Alan was jittered with a uniform distribution between  
190 500-1000 ms.

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

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

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

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

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

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

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

## 242 **Procedure**

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

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

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

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

## 270 **Scores**

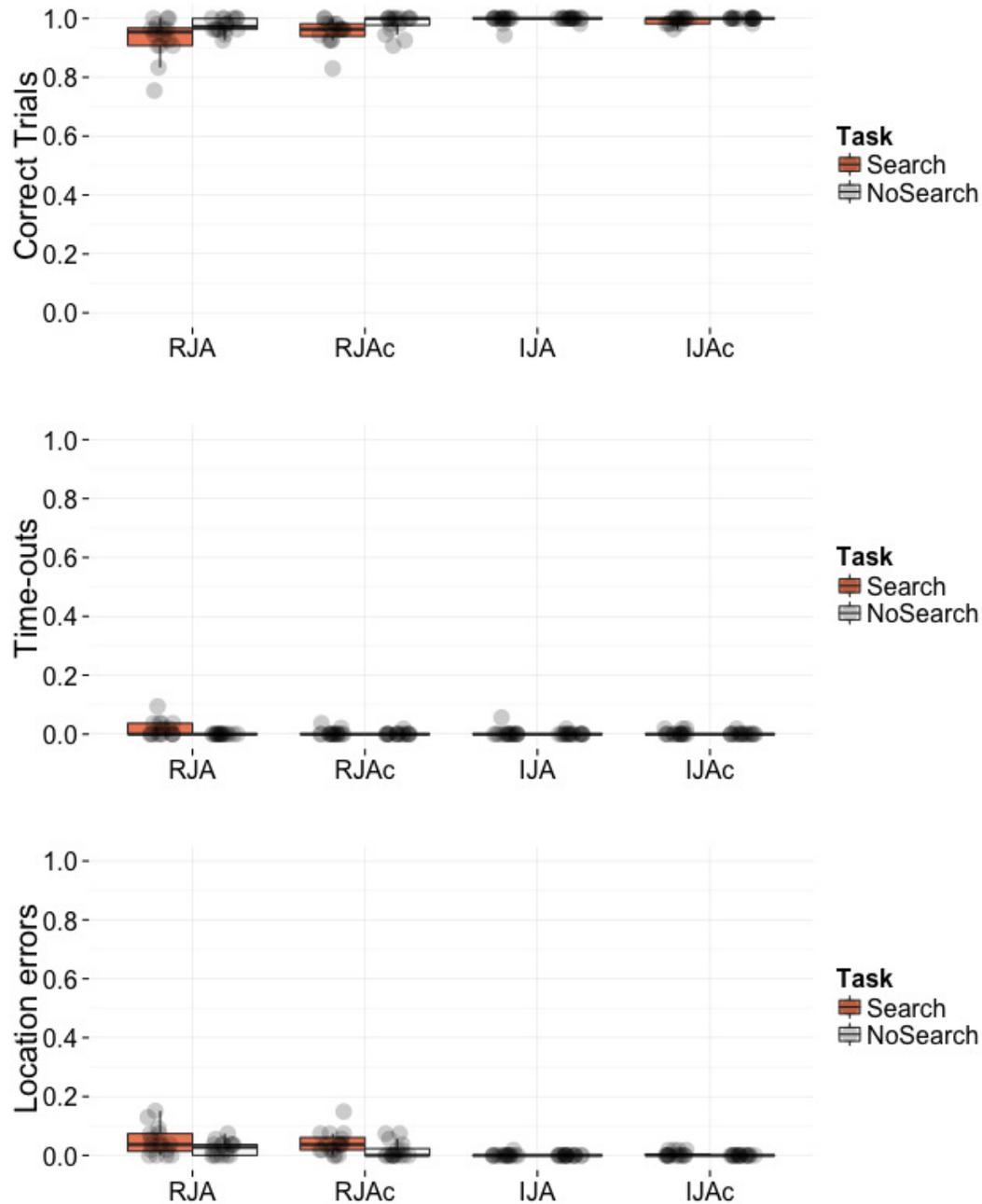
271           **Accuracy.** We calculated the proportion of trials where the participant succeeded in  
272 catching the burglar in each condition (i.e., RJA and RJAc) for each task separately (i.e., Search  
273 and NoSearch). We excluded from the accuracy analysis any trials that required a recalibration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

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

## 279 **Statistical Analyses**

280 Saccadic reaction times were analysed via repeated measures Analysis of Variance  
281 (ANOVA) using the ezANOVA (ez) package in R (Lawrence, 2013), reporting the generalised  
282 eta squared ( $\eta_G^2$ ) measure of effect size. Significant task\*condition interactions effects were  
283 followed-up with Welch's two sample unequal variances t-tests (Welch, 1947). As in our  
284 previous studies, we report analyses of the mean reaction time, having excluded trials with  
285 reaction times less than 150 ms as these are typically considered to be anticipatory responses.  
286 Trials timed out after 3000 ms providing a natural upper limit to reaction times. Full syntax and  
287 output for this analysis can be found in the Rmarkdown document (Supplementary Material 1).  
288 The RMarkdown also provides complementary ANOVAs of the mean and median of the  
289 untrimmed data, as well as a mixed random effects analysis using the lme4 R package (Bates,  
290 2005). The results of all analyses are consistent in terms of the predicted interaction between task  
291 and condition. A significance criterion of  $p < 0.05$  was used for all analyses.

293



294

295 *Figure 3.* Box plots displaying the proportion of correct trials, proportion of time-out errors, and proportion of

296 location errors. Data points represent individual participant means.

297

298

## Results

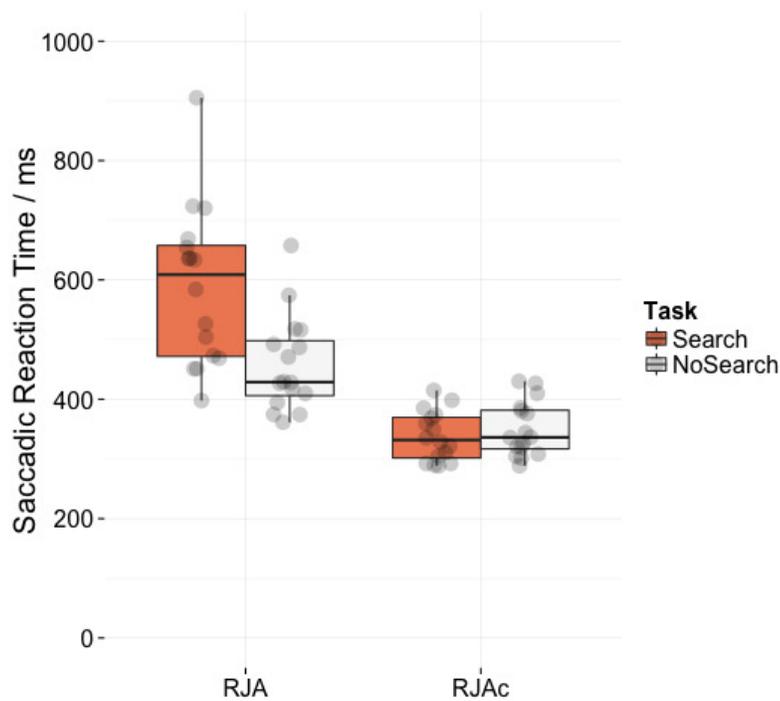
299 As depicted in Figure 3, participants performed at close-to-ceiling levels in terms of the  
300 trials successfully completed across all conditions. Of the small number of errors made, the  
301 majority were Location Errors in the RJA and RJAc conditions, whereby participants looked first  
302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

304 Figure 4 shows mean saccadic reaction times for correct trials in the RJA and RJAc  
305 conditions. Participants were significantly slower on the Search task than the NoSearch task  
306 (main effect of task ( $F(1,15) = 11.07, p = .005, \eta_G^2 = 0.13$ ). They were also significantly slower to  
307 respond on RJA trials than RJAc trials overall (main effect of condition,  $F(1,15) = 98.75, p <$   
308  $.001, \eta_G^2 = 0.57$ ). Importantly, there was also a significant task\*condition interaction ( $F(1,15) =$   
309  $43.86, p < .001, \eta_G^2 = 0.18$ ), indicating a larger effect of task in the RJA condition. This  
310 interaction was also present in all re-analyses of the data (see Supplementary Material 1).

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

318

319



*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.

320

**Discussion**

321

322

323

324

325

326

327

328

329

330

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

331

332

333

334

335

336

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

337

338

339

340

341

342

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

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

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

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

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

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

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

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

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

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

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

#### 414 **Summary**

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

426

**References**

- 427 Adamson, L. B., Bakeman, R., Deckner, D. F., & Romski, M. A. (2009). Joint engagement and  
428 the emergence of language in children with autism and Down syndrome. *Journal of*  
429 *Autism and Developmental Disorders*, 39(1), 84-96. doi: 10.1007/s10803-008-0601-7
- 430 Bayliss, A. P., di Pellegrino, G., & Tipper, S. P. (2005). Sex differences in eye gaze and  
431 symbolic cueing of attention. *Quarterly Journal of Experimental Psychology*, 58, 1-20.
- 432 Bates DM. 2005. Fitting linear mixed models in R. *R News* 5:27-30
- 433 Böckler, A., Timmermans, B., Sebanz, N., Vogeley, K., & Schilbach, L. (2014). Effects of  
434 observing eye contact on gaze following in high-functioning autism. *Journal of Autism*  
435 *and Developmental Disorders*, 44(7), 1651-1658. doi: 10.1007/s10803-014-2038-5
- 436 Bruinsma, Y., Koegel, R. L., & Koegel, L. K. (2004). Joint attention and children with autism: A  
437 review of the literature. *Mental Retardation & Developmental Disabilities Research*  
438 *Reviews*, 10(3), 169-175. doi: 10.1002/mrdd.20036
- 439 Bruner, J. S. (1974). From communication to language, A psychological perspective. *Cognition*,  
440 3(3), 255-287. doi: 10.1016/0010-0277(74)90012-2
- 441 Bruner, J. S. (1995). From joint attention to the meeting of minds: an introduction. In C. Moore  
442 & P. J. Dunham (Eds.), *Joint attention: its origins and role in development*. (pp. 1-14).  
443 Hillsdale, NJ: Lawrence Erlbaum Associates.
- 444 Caruana, N., Brock, J., & Woolgar, A. (2015). A frontotemporoparietal network common to  
445 initiating and responding to joint attention bids. *NeuroImage*, 108, 34-46. doi:  
446 10.1016/j.neuroimage.2014.12.041

- 447 Caruana, N., de Lissa, P., & McArthur, G. (2016). Beliefs about human agency influence the  
448 neural processing of gaze during joint attention. *Social Neuroscience*,  
449 doi:10.1080/17470919.2016.1160953
- 450 Caruana, N., Stieglitz Ham., H., Brock, J., Woolgar, A., Palermo, R., Kloth, N., & McArthur, G.  
451 (in press). Joint Attention Difficulties in Adults with Autism. *Autism*.
- 452 Cary, M. S. (1978). The role of gaze in the initiation of conversation. *Social Psychology*, 41(3),  
453 269-271. doi: 10.2307/3033565
- 454 Charman, T. (2003). Why is joint attention a pivotal skill in autism? *Philosophical Transactions*  
455 *Royal Society London Biological Sciences*, 358, 315-324. doi: 10.1098/rstb.2002.1199
- 456 Charman, T., Swettenham, J., Baron-Cohen, S., Cox, A., Baird, G., & Drew, A. (1997). Infants  
457 with autism: An investigation of empathy, pretend play, joint attention, and imitation.  
458 *Developmental Psychology*, 33(5), 781-789. doi: 10.1037/0012-1649.33.5.781
- 459 Danckert J, Ferber S, Doherty T, Steinmetz H, Nicolle D, Goodale MA (2002) Selective, non-  
460 lateralized impairment of motor imagery following right parietal damage. *Neurocase* 8,  
461 194–204
- 462 Dawson, G., Toth, K., Abbott, R., Osterling, J., Munson, J. A., Estes, A., & Liaw, J. (2004).  
463 Early social attention impairments in autism: Social orienting, joint attention, and  
464 attention to distress. *Developmental Psychology*, 40(2), 271-283. doi:  
465 <http://dx.doi.org/10.1037/0012-1649.40.2.271>
- 466 Frischen, A., Bayliss, A. P., & Tipper, S. P. (2007). Gaze cueing of attention: visual attention,  
467 social cognition, and individual differences. *Psychological Bulletin*, 133(4), 694-724. doi:  
468 10.1037/0033-2909.133.4.694

- 469 Gavrilov, Y., Rotem, S., Ofek, R., & Geva, R. (2012). Socio-cultural effects on children's  
470 initiation of joint attention. *Frontiers in Human Neuroscience*, *6*, 286.
- 471 Hamilton AFC, de, Grafton ST. 2008. Action outcomes are represented in human inferior  
472 frontoparietal cortex. *Cerebral Cortex*, *18* (5), 1160–1168.
- 473 Lawrence MA (2013). ez: Easy analysis and visualization of factorial experiments. R package  
474 version 4.2-2, URL <http://CRAN.R-project.org/package=ez>.
- 475 Leekam, S. (2015). Social cognitive impairment and autism: what are we trying to explain?  
476 *Philosophical Transactions of the Royal Society of London B: Biological Sciences*, *371*,  
477 1686.
- 478 Lord, C., Risi, S., Lambrecht, L., Cook, E.H., Leventhal, B.L., DiLavore, P.C., Pickles, A.,  
479 Rutter, M., 2000. The autism diagnostic observation schedule—generic: a standard  
480 measure of social and communication deficits associated with the spectrum of autism.  
481 *Journal of Autism and Developmental Disorders*, *30*(3), 205–223. doi:  
482 10.1023/A:1005592401947
- 483 Loveland, K. A., & Landry, S. H. (1986). Joint attention and language in autism and  
484 developmental language delay. *Journal of Autism and Developmental Disorders*, *16*(3),  
485 335-349. doi: 10.1007/BF01531663
- 486 Mundy, P., Sigman, M., & Kasari, C. (1990). A longitudinal study of joint attention and  
487 language development in autistic children. *Journal of Autism and Developmental*  
488 *Disorders*, *20*(1), 115-128. doi: 10.1007/bf02206861
- 489 Mundy, P., Block, J., Vaughan Van Hecke, A., Delgado, C., Parlade, M., Pomeroy, Y. (2007).  
490 Individual differences in the development of joint attention in infancy. *Child*  
491 *Development*, *78*, 938–954.

- 492 Murray, D. S., Creaghead, N. A., Manning-Courtney, P., Shear, P. K., Bean, J., & Prendeville, J.  
493 A. (2008). The relationship between joint attention and language in children with autism  
494 spectrum disorders. *Focus on Autism & Other Developmental Disabilities*, 23(1), 5-14.  
495 doi: 10.1177/1088357607311443
- 496 Nation, K., & Penny, S. (2008). Sensitivity to eye gaze in autism: Is it normal? Is it automatic? Is  
497 it social? *Development and Psychopathology*, 20(01), 79-97. doi:  
498 doi:10.1017/S0954579408000047
- 499 Oberwelling, E., Schilbach, L., Barisic, I., Krall, S. C., Vogeley, K., Fink, G. R., Herpertz-  
500 Dahlmann, B., Konrad, K., Schulte-Rüther, M. (2016). Look into my eyes: Investigating  
501 joint attention using interactive eye-tracking and fMRI in a developmental sample.  
502 *NeuroImage*, in press.
- 503 Osterling, J., Dawson, G., & Munson, J. A. (2002). Early recognition of 1-year-old infants with  
504 autism spectrum disorder versus mental retardation. *Development and Psychopathology*,  
505 14(02), 239-251. doi: doi:10.1017/S0954579402002031
- 506 Redcay, E., Dodell-Feder, D., Mavros, P. L., Kleiner, A. M., Pearrow, M. J., Triantafyllou, C.,  
507 Gabrieli, J. D., & Saxe, R. (2012). Atypical brain activation patterns during a face-to-face  
508 joint attention game in adults with autism spectrum disorder. *Human Brain Mapping*, 34,  
509 2511-2523. doi: 10.1002/hbm.22086.
- 510 Saito, D. N., Tanabe, H. C., Izuma, K., Hayashi, M. J., Morito, Y., Komeda, H., . . . Sadato, N.  
511 (2010). "Stay Tuned": inter-individual neural synchronization during mutual gaze and  
512 joint attention. *Frontiers in Integrative Neuroscience*, 4, 127. doi:  
513 10.3389/fnint.2010.00127

- 514 Saxe, R., & Kanwisher, N. (2003). People thinking about thinking people: The role of the  
515 temporo-parietal junction in “theory of mind”. *NeuroImage*, *19*(4), 1835–1842. doi:  
516 [http://dx.doi.org/10.1016/S1053-8119\(03\)00230-1](http://dx.doi.org/10.1016/S1053-8119(03)00230-1).
- 517 Saxon, Terrill F., and John T. Reilly. (1999). Joint Attention and Toddler Characteristics: Race,  
518 Sex and Socioeconomic Status. *Early Child Development and Care*, *149*(1), 59-69.
- 519 Schilbach, L., Wilms, M., Eickhoff, S. B., Romanzetti, S., Tepest, R., Bente, G., Shah N. J.,  
520 Fink, G. R., & Vogeley, K. (2010). Minds made for sharing: Initiating joint attention  
521 recruits reward-related neurocircuitry. *Journal of Cognitive Neuroscience*, *22*(12), 2702-  
522 2715. doi: 10.1162/jocn.2009.21401.
- 523 Schilbach, L., Timmermans, B., Reddy, V., Costall, A., Bente, G., Schlicht, T., & Vogeley, K.  
524 (2013). Toward a second-person neuroscience. *Behavioral and Brain Sciences*, *36*(4),  
525 393-414. doi: 10.1017/S0140525X12000660.
- 526 Senju, A., & Johnson, M. H. (2009). The eye contact effect: mechanisms and development.  
527 *Trends in Cognitive Sciences*, *13*(3), 127-134. doi:  
528 <http://dx.doi.org/10.1016/j.tics.2008.11.009>
- 529 Singular Inversions. (2008). FaceGen Modeller (Version 3.3) [Computer Software]. Toronto,  
530 ON: Singular Inversions.
- 531 SR Research. (2004). Experiment Builder (Version 1.10.165). Ontario.
- 532 Stone, W. L., Ousley, O. Y., & Littleford, C. D. (1997). Motor imitation in young children with  
533 autism: what's the object? *Journal of Abnormal Child Psychology*, *25*, 475-485. doi:  
534 10.1023/A:1022685731726

- 535 Tomasello, M. (1995). Joint attention as social cognition. In C. Moore & P. J. Dunham (Eds.),  
536 *Joint Attention: Its Origins and Role in Development*. Hillsdale: Lawrence Erlbaum  
537 Associates.
- 538 Welch, B. L. (1947). The generalization of 'student's' problem when several different population  
539 variances are involved. *Biometrika*, 34, 28-35.
- 540 Wong, C., & Kasari, C. (2012). Play and joint attention of children with autism in the preschool  
541 special education classroom. *Journal of Autism and Developmental Disorders*, 42(10),  
542 2152-2161. doi: 10.1007/s10803-012-1467-2
- 543 Wykowska, A., Wiese, E., Prosser, A., Müller, H. J., & Hamed, S. B. (2014). Beliefs about the  
544 minds of others influence how we process sensory information. *PLoS ONE*, 9(4), e94339.  
545 doi:10.1371/journal.pone.0094339