

---

*Research Article: New Research | Cognition and Behavior*

## **Frontal, parietal and temporal brain areas are differentially activated when disambiguating potential objects of joint attention**

<https://doi.org/10.1523/ENEURO.0437-19.2020>

**Cite as:** eNeuro 2020; 10.1523/ENEURO.0437-19.2020

Received: 22 October 2019

Revised: 7 July 2020

Accepted: 4 August 2020

---

*This Early Release article has been peer-reviewed and accepted, but has not been through the composition and copyediting processes. The final version may differ slightly in style or formatting and will contain links to any extended data.*

**Alerts:** Sign up at [www.eneuro.org/alerts](http://www.eneuro.org/alerts) to receive customized email alerts when the fully formatted version of this article is published.

Copyright © 2020 Kraemer et al.

This is an open-access article distributed under the terms of the Creative Commons Attribution 4.0 International license, which permits unrestricted use, distribution and reproduction in any medium provided that the original work is properly attributed.

1 1. Frontal, parietal and temporal brain areas are differentially activated when disambiguating  
2 potential objects of joint attention

3

4 2. Neural Substrates of Joint Attention

5 3. P. M. Kraemer<sup>1,5,†</sup>, M. Görner<sup>1,2,3,†</sup>, H. Ramezanpour<sup>1,2,3,†</sup>, P. W. Dicke<sup>1</sup> and P. Thier<sup>1,4,\*</sup>

6 <sup>1</sup>Department of Cognitive Neurology, Hertie Institute for Clinical Brain Research, University of Tübingen, 72076  
7 Tübingen, Germany.

8 <sup>2</sup>Graduate School of Neural and Behavioural Sciences, University of Tübingen, 72074 Tübingen, Germany.

9 <sup>3</sup>International Max Planck Research School for Cognitive and Systems Neuroscience, University of Tübingen,  
10 72074 Tübingen, Germany.

11 <sup>4</sup>Werner Reichardt Centre for Integrative Neuroscience, University of Tübingen, 72076 Tübingen, Germany.

12 <sup>5</sup>Center for Decision Neuroscience, Faculty of Psychology, University of Basel, 4055 Basel, Switzerland.

13

14 4. † These authors contributed equally to this work. PT developed the conceptual framework of  
15 the research. PT, PK and HR designed the experiments. PK and PWD performed the  
16 experiments. PK and MG analyzed the data. All authors contributed to the interpretation of  
17 results and the writing.

18 5. Corresponding Author: Peter Thier, Department of Cognitive Neurology, Hertie Institute for  
19 Clinical Brain Research, Hoppe-Seyley-Str. 3, 72076 Tübingen, Germany. E-mail: [thier@uni-](mailto:thier@uni-tuebingen.de)  
20 [Tübingen.de](mailto:thier@uni-tuebingen.de).

21 6. Number of main figures: 5

22 7. Extended data: 6 Figures, 0 Tables, 1 Paradigm description

23 8. Number of words of Abstract: 249

24 10. Number of words for Significance Statement: 108

25 11. Number of words for Introduction: 607

26 12. Number of words for Discussion: 1830

27 13. Acknowledgments: We are grateful to Friedemann Bunjes and Michael Erb for technical  
28 support.

29 14. Conflicts of interest: Authors report no conflict of interest

30 15. Funding sources: This work was supported by the Deutsche Forschungsgemeinschaft (TH  
31 425/12-2).

32

33 **Abstract**

34 Humans establish joint attention with others by following the other's gaze. Previous work has  
35 suggested that a cortical patch ("gaze following patch," GFP) close to the posterior superior  
36 temporal sulcus may serve as a link between the extraction of the other's gaze direction and the  
37 resulting shifts of attention, mediated by human LIP (hLIP). However, it is not clear how the  
38 brain copes with situations in which information on gaze direction alone is insufficient to  
39 identify the target object because more than one may lie along the gaze vector. In this fMRI  
40 study we tested human subjects on a paradigm that allowed the identification of a target object  
41 based on the integration of the other's gaze direction and information provided by an auditory  
42 cue on the relevant object category. Whereas the GFP activity turned out to be fully determined  
43 by the use of gaze direction, activity in hLIP reflected the total information needed to pinpoint  
44 the target. Moreover, in an exploratory analysis, we found that a region in the inferior frontal  
45 junction (IFJ) was sensitive to the total information on the target. An examination of the BOLD  
46 time courses in the three identified areas suggests functionally complementary roles. Whereas  
47 the GFP seems to primarily process directional information stemming from the other's gaze, the  
48 IFJ may help to analyze the scene when gaze direction and auditory information are not  
49 sufficient to pinpoint the target. Finally, hLIP integrates both streams of information in order to  
50 shift attention to distinct spatial locations.

51

52

53

54

55 **Significance statement**

56 Our paper captures work deploying fMRI to identify brain structures and mechanisms that allow  
57 us to pinpoint the object attended by the other out of the many hit by her/his gaze vector. Our  
58 results suggest that an area at the border between premotor and prefrontal cortex plays a major  
59 role in providing the complementary information needed by the temporo-parietal underpinnings  
60 of gaze following to shift the observer's attention to the correct object. Our results support the  
61 interplay of a network of distributed elements with distinct functional contributions allowing us  
62 to deploy *Joint Attention* – a key underpinning of viable social behavior and a full-fledged  
63 *Theory of the (other's) Mind*.

64

65 **Code Accessibility**

66 All data and analysis scripts are freely available on  
67 [https://figshare.com/projects/Contextual\\_Gaze\\_Following/78222](https://figshare.com/projects/Contextual_Gaze_Following/78222)

68

69 **Introduction**

70 We follow the other's gaze to objects of her/his attention which allows us to shift our own  
71 attention to the same object, and to thereby establish joint attention. By associating our object-  
72 related intentions, expectations and desires with the other one, joint attention allows us to  
73 develop a *Theory of (the other's) Mind* (TOM) (Emery, 2000). TOM is a major basis of  
74 successful social interactions (Baron-Cohen, 1995, 1994) and, arguably, its absence is at the core  
75 of neuropsychiatric disorders such as autism.

76 Human gaze following is geometric (Atabaki et al., 2015; Butterworth and Jarrett, 1991). This  
77 means that we use the other's gaze vector to identify the exact location of the object of interest.  
78 The features of the human eye such as the high contrast between the white sclera and dark iris  
79 allow us to determine the other's eye direction at high resolution (Bock et al., 2008; Kobayashi  
80 and Kohshima, 1997). However, knowledge of direction is not sufficient to pinpoint an object in  
81 3D. In principle, differences between the directions of the two eyes, i.e. knowledge of the  
82 vergence angle, could be exploited to this end. Yet, this will work only for objects close to the  
83 beholder as the angle will become imperceptibly small if the objects are outside the confines of  
84 peripersonal space. On the other hand, gaze following remains precise also for objects quite far  
85 from the other one although the gaze vector will in many cases hit more than one object  
86 (Butterworth and Jarrett, 1991). Hence, how can these objects be disambiguated? We  
87 hypothesized that singling out the relevant object is a consequence of recourse to prior  
88 information on the objects and their potential value for the other. For instance, let us assume that  
89 the day is hot and that the other's appearance may suggest thirst and the desire to take a sip of  
90 something cool. If her/his gaze hits a cool beverage within a set of other objects of little

91 relevance for a thirsty person, the observer might safely infer that the beverage is the object of  
92 desire. In this example, gaze following is dependent on prior assumptions about the value of  
93 objects for the other. Of course, the value the object may have for the observer also matters. For  
94 instance, Liuzza et al. showed that an observer's appetite to follow the other's gaze to portraits  
95 of political leaders is modulated by the degree of political closeness (Liuzza et al., 2011). If the  
96 politician attended by the other was a political opponent of the observer, the willingness to  
97 follow gaze was significantly reduced. Also knowing that gaze following may be inadequate in a  
98 given situation and that the other may become aware of an inadequate behavior will suppress it  
99 (Teufel et al., 2010, 2009). However, only assumptions about the object value of the other will  
100 help to disambiguate the scene.

101 Following the gaze of others to a particular object is accompanied by a selective BOLD signal in  
102 an island of cortex in the posterior superior temporal sulcus (pSTS), the "gaze-following patch"  
103 (GFP) (Laube et al., 2011; Marquardt et al., 2017; Materna et al., 2008). In these studies, the  
104 target object could be identified unambiguously by gaze direction as for a given gaze direction  
105 the vector hit one object only. Hence, it remains unclear if the GFP helps to integrate the  
106 information needed to disambiguate the object choice in case the gaze vector hits more than one  
107 object. In order to address this question, we carried out an fMRI study in which the selection of  
108 the object of joint attention required that the observer resorted to another source of information  
109 aside from the gaze cue.

110

111 **Materials and Methods**

112 *Participants*

113 Nineteen healthy, right-handed volunteers (9 females and 10 males, mean age 27.4, SD = 3.6)  
114 participated in the study over three sessions. Participants gave written consent to the procedures  
115 of the experiment. The study was approved by the local Ethics Review Board and was carried  
116 out in accordance with the principles of human research ethics of the Declaration of Helsinki.

117

118 *Task and procedure*

119 The study was conducted in three sessions across separate days. On day 1, we instructed  
120 participants about the study goals and familiarized them with the experimental paradigms outside  
121 the MRI-scanner by carrying out all relevant parts of the fMRI experiments. The following  
122 fMRI-experiments included a functional localizer paradigm for the scanning session on day 2 as  
123 well as a contextual gaze following paradigm for the scanning session on day 3.

124 Behavioral session. After participants had been familiarized with the tasks, they were head-fixed  
125 using a chinrest and a strap to fix the forehead to the rest. Subjects were facing towards a  
126 frontoparallel screen (resolution = 1280×1024 pixels, 60 Hz) (distance to eyes ≈ 600 mm). Eye  
127 tracking data were recorded while participants had to complete 80 trials of the localizer paradigm  
128 and 72 trials of contextual gaze following.

129 Localizer task. We resorted to the same paradigm used in (Marquardt et al., 2017) to localize the  
130 gaze following network and in particular its core, the GFP. In this paradigm, subjects were asked  
131 to make saccades to distinct spatial targets based on information provided by a human portrait  
132 presented to the observer. Depending on the instruction, subjects either had to rely on the seen  
133 gaze direction to identify the correct target (*gaze following* condition) or, alternatively, they had

134 to use the color of the irises, changing from trial to trial but always mapping to one of the targets,  
135 in order to make a saccade to the target having the same color (*color mapping* condition). In  
136 other words, the only difference between the two tasks was the information subjects had to  
137 exploit in order to solve the task, while the visual stimuli were the same.

138 This task is associated with higher BOLD activity in the GFP, a region close to the posterior end  
139 of the superior temporal sulcus, when subjects performed gaze following compared to color  
140 mapping. The task is further associated with the activation of regions in the posterior parietal  
141 cortex as well as the frontal cortex that take part in controlling spatial attention and saccade  
142 generation (Marquardt et al., 2017; Materna et al., 2008). Out of the 19 subjects of our study, 16  
143 performed 6 runs (40 trials per run) and for reasons of time management during image  
144 acquisition, one subject performed 5 runs and two subjects performed 4 runs.

145 Contextual gaze following task. An example of a trial is shown in Fig. 1. Each trial consisted of  
146 the following sequence of events. The trial started with the appearance of the portrait of an avatar  
147 (6.7×10.5 degrees of visual angle) image in the center of the screen together with four arrays of  
148 drawn objects (houses and hands, 3 objects per array). Subjects were asked to fixate on a red  
149 fixation dot (diameter) between the avatar's eyes. After 5 seconds of baseline fixation, the  
150 avatar's eye gaze shifted towards one specific target object. Simultaneously, a spoken instruction  
151 either specified the object class of the target (spoken words "hand" or "house") or was not  
152 informative ("none"). While maintaining fixation, subjects needed to judge which object the  
153 avatar was looking at. After 5 seconds delay, the fixation dot vanished, an event that served as  
154 the *go*-signal. Participants had 2 seconds to make a saccade to the chosen target object and fixate  
155 it until a subsequent blank fixation screen was presented for about 8 seconds. The subjects were  
156 instructed to perform the task as accurately as possible. They were specifically instructed, when

157 unsure about the actual target, nevertheless to rely on gaze and contextual information and  
158 choose the target they believed the avatar to be looking at.

159 The information provided by the spoken instruction distinguished three experimental conditions,  
160 one *unambiguous*, and two ambiguous conditions: *ambiguous-informative* and *ambiguous-*  
161 *uninformative*. The verbal instruction in the *unambiguous* condition reduced the number of  
162 potential targets from three to one by naming the object category with only one representative in  
163 the array. For instance, in Fig. 1, the avatar gazes at the lower left array, specifying two hands  
164 and one house as potential gaze targets. An unambiguous instruction would be the auditory cue  
165 “house”. The *ambiguous-informative* instruction in this example, “hand”, reduced the number of  
166 potential gaze targets to two. In the *ambiguous-uninformative* condition, the instruction would  
167 have been “none”, not suited to reduce the number of potentials targets.

168 Participants performed six blocks of 30 trials each (10 per condition), summing up to 180 trials  
169 in total.

170

### 171 *Stimuli*

172 Visual and auditory stimuli as well as data collection was controlled by a custom-made Linux  
173 based system. The stimuli in the localizer task were designed according to the stimuli used in a  
174 previous study (Marquardt et al., 2017). The stimuli of the contextual gaze following task  
175 consisted of an avatar and in total 12 target objects from two categories (houses and hands). The  
176 avatar was generated with a custom-made OpenGL library which offers a controlled virtual 3D-  
177 environment in which an avatar can be set to precisely gaze at specific objects. More  
178 specifically, the program allows objects to be placed on a circle, parallel to the coronal axis, and  
179 anterior to the avatar face. For each stimulus, we placed 12 objects in the surroundings of the

180 avatar. The location of individual objects was fully determined by the distance to the coronal  
181 plane at the level of the avatar's nasion, the radius of the circle and the angle of the object on that  
182 circle. By keeping the angle on the circle constant for each set of three objects, we created four  
183 arrays at angles 120°, 150°, 210° and 240°. The individual locations of these objects were  
184 specified by varying the distance and the circle radii based on trigonometric calculations. For  
185 these calculations we assumed a right triangle from the avatar's nasion with the hypotenuse  
186 pointing towards the object, an adjacent leg (length corresponded to the distance of the circle)  
187 proceeding orthogonal to the coronal plane, and an opposite leg which corresponded to the  
188 radius. By keeping  $\tan\alpha$  fixed to 0.268, we varied the distances and circle radii. For the 120° and  
189 240° arrays, the circle radii were 335, 480, 580 and the distances were 90, 129 and 151 virtual  
190 mm. For the 150° and 210° arrays, the radii were 380, 510 and 590 and the distances were 102,  
191 137 and 158 virtual mm. The reason for the difference of radii and distances between 120°/240°  
192 and 150°/210° arrays was that this allowed the total width of the screen to be exploited. This  
193 procedure guaranteed that the angle of the gaze vector to all objects on the array was almost  
194 identical. This makes it relevant to take contextual information into account in order to choose  
195 the true target.

196 The objects were drawings of the two categories houses and hands, downloaded from freely  
197 available online sources (<http://www.allvectors.com/house-vector/>,  
198 [https://www.freepik.com/free-vector/hand-drawn-  
199 hands\\_812824.htm#term=hands&page=1&%20position=37](https://www.freepik.com/free-vector/hand-drawn-hands_812824.htm#term=hands&page=1&%20position=37)). The target objects were arranged in  
200 four radial directions (three objects in each direction) with the avatar's eyes as the origin; in  
201 other words, the avatar's gaze always hit one out of three objects along the gaze vector though  
202 participants were not able to tell which of the three it was. On each array, either 2 hands and one

203 house or one hand and two houses were present. Further, we fixed the number of hands and  
204 houses per hemifield to three. The relative order of the objects was pseudo-randomized from trial  
205 to trial.

206 We created a pool of stimulus sets which satisfied three constraints: There was an equal number  
207 of trials in which a) the targets were hands or houses, b) targets were presented with an  
208 *unambiguous*, *ambiguous-informative* and *ambiguous-uninformative* instruction, and c) the  
209 spatial position (one out of twelve potential positions) of targets was matched. This led to 72  
210 stimulus sets. We exposed every subject to 180 trials in which each stimulus set was shown  
211 twice and for the residual 36 trials, stimuli were drawn pseudo-randomly from the stimulus pool  
212 so that the three aforementioned criteria were met.

213 Auditory instructions were delivered via headphones (Sennheiser HD 201, Wedemark-  
214 Wennebostel, Germany during the behavioral session, and standard air pressure headphones of  
215 the scanner system during the MRI sessions). The auditory instructions “hand”, “house” and  
216 “none” were computer generated with the web application imTranslator  
217 (<http://imtranslator.net/translate-and-speak/speak/english/>) and processed with the software  
218 Audacity 2.1.2. The sound files had a duration of 600 ms.

219

#### 220 *Eye tracking*

221 During all three sessions, we recorded eye movements of the right eyes using commercial eye  
222 tracking systems (Behavioral sessions: Chronos Vision C-ETD, Berlin, Germany, sampling rate  
223 400 Hz, resolution  $< 1^\circ$  visual angle; Scanning sessions: SMI iView X MRI-LR, Berlin,  
224 Germany, sampling rate = 50 Hz, resolution  $\approx 1^\circ$  visual angle).

225 Eye tracking data was processed as follows. First, we normalized the raw eye tracking signal by  
226 dividing it by the average of the time series. Eye blinks were removed using a velocity threshold  
227 ( $> 1000$  %/s visual angle). Next, we focused on a time window in which we expected the saccades  
228 to the target objects to occur ( $[go\text{-signal} - 500 \text{ ms}, go\text{-signal} + 1800 \text{ ms}]$ ). Within this time  
229 window, we detected saccades by identifying the time point of maximal eye movement velocity.  
230 Pre- and post-saccadic fixation positions were determined by averaging periods of 200 ms before  
231 and after the saccade occurred. Due partly to extensive noise of the eye tracking signal recorded  
232 in the scanner, we did not automatize the categorization of the final gaze position. Instead, we  
233 plotted X- and Y coordinates of the post-saccadic eye position for every run that was not  
234 contaminated by too much noise. An investigator (MG), who was blind to the true gaze target-  
235 directions of the stimulus face, manually validated which trials yielded positions that were  
236 clearly assignable to a distinct object location. For the behavioral analysis we only used valid  
237 trials (mean number of valid trials per participant = 80.2,  $SD = 45.4$ , range = [0,153]) and  
238 weighted the individual performance values by its number in order to compute weighted means  
239 and  $SDs$ . Note that we used these valid trials only for the behavioral analysis but used all trials of  
240 the participants for the fMRI analysis, assuming that eye tracking measurement noise was  
241 independent of the performance of the subjects.

242

243 *fMRI acquisition and preprocessing.*

244 We acquired MR images using a 3T scanner (Siemens Magnetom Prisma, Erlangen, Germany)  
245 with a 20-channel phased array head coil. The head of the subjects was fixed inside the head coil  
246 by using plastic foam cushions to avoid head movements. An AutoAlign sequence was used to  
247 standardize the alignment of images across sessions and subjects. A high-resolution T1-weighted

248 anatomical scan (MP-RAGE, 176×256×256 voxel, voxel size 1×1×1 mm) and local field maps  
249 were acquired. Functional scans were carried out using a T2<sup>\*</sup>-weighted echo-planar multi-banded  
250 2D sequence (multi-band factor = 2, TE = 35 ms, TR = 1500 ms, flip angle = 70°) which covered  
251 the whole brain (44×64×64 voxel, voxel size 3×3×3 mm, interleaved slice acquisition, no gap).  
252 For image preprocessing we used the MATLAB SPM12 toolbox (Statistical Parametric  
253 Mapping, <https://www.fil.ion.ucl.ac.uk/spm/>). The anatomical images were segmented and  
254 realigned to the SPM T1 template in MNI space. The functional images were realigned to the  
255 first image of each respective run, slice-time corrected and coregistered to the anatomical image.  
256 Structural and functional images were spatially normalized to MNI space. Finally, functional  
257 images were spatially smoothed with a Gaussian kernel (6 mm full-width at half maximum).

258

259 *fMRI analysis.*

260 We estimated a generalized linear model (GLM) to identify regions of interest (ROIs) of single  
261 subjects. On these regions, we performed time course analyses to investigate event-related  
262 BOLD signal changes. In a first-level analysis, we estimated GLMs for the localizer task  
263 (GLM<sub>loc</sub>) and the contextual gaze following task (GLM<sub>cgf</sub>). The GLM<sub>loc</sub> included predictors for  
264 the onset of directional cues and of the baseline fixation phase. The GLM<sub>cgf</sub> had predictors for  
265 the onset of the contextual instruction coinciding with the gaze cue. These event specific  
266 predictors of the two GLMs used the canonical hemodynamic response function of SPM to  
267 model the data. We corrected for head motion artifacts by the estimation of six movement  
268 parameters with the data of the realignment preprocessing step. Low-frequency drifts were  
269 filtered using a high-pass filter (cutoff at 1/128 Hz).

270 GFP and hLIP localizer. Before collecting the data, we specified the expected locations of two  
271 brain areas from the literature. We resorted to the parietal coordinates of the human homologue  
272 of the monkey area LIP (hLIP) which had been identified using a delayed saccade task (Serenio et  
273 al., 2001). The GFP standard coordinates were taken from Marquardt, Ramezani et al.  
274 (2017). We transformed the standard coordinates for the hLIP and the GFP from Talairach space  
275 into MNI space, using an online transformation method of Lacadie and colleagues (Lacadie et  
276 al., 2008) (<http://sprout022.sprout.yale.edu/mni2tal/mni2tal.html>).

277 To identify ROIs at the group level, we compared beta weights of the statistical parametric maps  
278 from the  $GLM_{loc}$  in a second-level analysis. The GFP weights were derived from the contrast  
279 *gaze following vs. color mapping*, and the hLIP weights from the contrast *gaze following vs.*  
280 *baseline fixation*. To be characterized as GFP or hLIP, a cluster's maximum weight had to be  
281 located in close proximity to their respective standard coordinates.

282 We aimed to identify ROIs on an individual subject level. To this end, we used the contrast maps  
283 from the first-level analysis of the  $GLM_{loc}$ . We selected the coordinates of the maximum contrast  
284 voxel which minimized the distance to the group level coordinates. This voxel had to be part of a  
285 statistically significant cluster (cluster size  $\geq 6$ ,  $p < 0.05$ ). Due to relatively low signal-to-noise  
286 ratio in the *gaze following vs. color mapping* contrast, and the corresponding increased risk of  
287 false-positive activations, we decided to introduce a second criterion to make GFP localization  
288 more rigorous in single subjects. This *proximity criterion* additionally required the cluster to be  
289 located at least partially within 10 mm from the group level coordinates of the respective ROI.

290 Contrasts of context conditions. In addition to our a priori ROIs, we were interested in whether  
291 the contextual gaze following task might also activate regions which we did not consider  
292 beforehand. We performed a whole-brain analysis on the data from the contextual gaze following

293 task. Using the  $GLM_{\text{cgl}}$ , we contrasted the weights of the two *ambiguous* conditions with the  
294 *unambiguous* condition at the group level (second-level analysis, significance threshold  $p <$   
295  $0.001$ , cluster size  $\geq 6$  voxel) as well as at the single subject level (first-level analysis,  
296 significance threshold  $p < .05$ , cluster size  $\geq 6$  voxel).

297 Time course analysis. We determined the individual time courses of the BOLD signal within  
298 sphere-shaped ROIs. Whenever we identified an ROI on the single-subject level, spheres with a  
299 radius of 5 mm were centered at the individual ROI coordinates. In case the identification of a  
300 ROI on the single-subject level was not possible, we deployed spheres with a radius of 10 mm  
301 centered at the group level location, assuming these spheres would capture relevant single-  
302 subject activity.

303 For every subject and sphere, raw time series of the BOLD signal were extracted using the  
304 MATLAB toolbox *marsbar 0.44* (<http://marsbar.sourceforge.net>). Due to technical problems in  
305 the reconstruction of trial times, for five participants we included only five runs and for two only  
306 four runs into the analysis. The time course of every trial was normalized by the average signal  
307 intensity 5 s before the onset of the contextual instruction and transformed into % of signal  
308 change. For each participant, we averaged time courses across trials and runs and used the time  
309 courses of the three contextual conditions in the six ROIs for our analysis. To test differences  
310 across conditions for statistical significance, we performed permutation tests at each time point  
311 after contextual instruction delivery. To do so, we pooled the data of two experimental  
312 conditions, respectively, and produced 10,000 random splits for each pool. By computing the  
313 differences between the means of these splits, we obtained a distribution of differences under the  
314 null hypothesis. Calculating the fraction of values more extreme than the actual difference  
315 between means allowed us to obtain a  $p$ -value for each time bin. To account for the multiple

316 comparison problem, we transformed  $p$ -values into FDR corrected  $q$ -values (Benjamini and  
317 Hochberg, 1995) and considered each time bin with  $q < .05$  as statistically significant.

318 We carried out two additional analyses, the first one to obtain Bayesian credible intervals (BCI)  
319 for the time courses and a second one, a “searchlight” analysis which scanned the whole brain  
320 for voxels whose signal could be used to decode the *unambiguous* vs. the *ambiguous-*  
321 *uninformative* conditions.

322 For the first we estimated hierarchical models for each experimental condition and ROI allowing  
323 the intercept to vary for each participant. The models were linear combinations of 7 sinusoidal  
324 basis functions. Model estimation was conducted using the *nideconv* package (de Hollander and  
325 Knäpen 2017) which interfaces with the Stan probabilistic programming language for Bayesian  
326 model estimation (Stan Development Team, 2018). 95% BCI indicate that, given the data and the  
327 model, the data generating parameter is included in the interval with a probability of 95%.

328 For the second we used the beta-images from the first-level analysis as input to a searchlight  
329 algorithm that estimated classification accuracies for each voxel. We employed the TDT toolbox  
330 (Hebart et al., 2015) and trained SVMs (leave-one-run-out cross-validated) to classify  
331 *unambiguous* and *ambiguous-uninformative* conditions. This yielded one map per participant  
332 representing the estimated classification accuracies minus the chance-level of 0.5. For a  
333 subsequent second-level analysis these individual maps were then smoothed (4mm kernel) and  
334 fed into SPM’s second-level pipeline to obtain a group level t-map as suggested by Hebart et al.,  
335 2015. Only voxels with a value corresponding to a  $p$ -value smaller or equal to 0.001 were  
336 included into the end result. Since classification accuracies do not follow a normal distribution  
337 computing a t-map is not the ideal method as it is likely to provide an overestimation of  
338 significance. However, even if an approach tending to overestimate statistical significances fails

339 to detect areas of significant classification, this method can still safely be considered  
340 conservative. This was the case with respect to GFP, whose possible involvement in the  
341 differentiation of the various conditions we wanted to reexamine using the searchlight analysis.  
342 As a final step we used the same ROIs that were used in the time course analysis to extract the  
343 local accuracies (minus chance-level) from each participant's accuracy map. The obtained  
344 distributions across participants were tested against the result of an ideal ignorant classifier that  
345 always performs at chance level using a Wilcoxon signed-rank test.

346

347

348 **Results**

349

350 Our subjects participated in two fMRI experiments. The first one was a *localizer* task that  
351 allowed us to identify two regions of interest of which we know are relevant for attentional shifts  
352 based on social cues, the GFP and hLIP (Marquardt et al., 2017; Materna et al., 2008). Our main  
353 intention was to investigate the BOLD activity of these regions in a *contextual gaze following*  
354 task (experiment 2). In this task, the subjects used the gaze direction of a human avatar,  
355 complemented by a spoken instruction. In one out of three conditions the observer was able to  
356 unambiguously identify the relevant object out of several hit by the other's gaze vector. This was  
357 the case in the *unambiguous* condition in which the spoken instruction identified an object class  
358 represented by only one exemplar on the avatar's gaze vector. In the two other conditions (to  
359 which we refer collectively as ambiguous conditions) the spoken information was insufficient  
360 either because two exemplars of the relevant object category were available (*ambiguous-*  
361 *informative* condition) or because the verbal instruction was uninformative (*ambiguous-*  
362 *uninformative* condition). In the latter case, observers were left with a choice between three  
363 objects.

364

365 *Behavioral Performance.* In the localizer task, subjects were able to hit targets reliably and  
366 without significant difference between the two conditions (median hit rates: *gaze following*:  $0.94$   
367  $\pm 0.13$  *SD*; *color matching*:  $0.92 \pm 0.09$  *SD*;  $p = 0.6$ , two-tailed t-test,  $N = 19$ , Fig. 2). Using the  
368 gaze following performance in the localizer task as reference we assumed the following expected  
369 hit rates for the contextual gaze following task:  $0.94$  for the *unambiguous condition*,  $0.94 * 1/2$  for  
370 the *ambiguous-informative* and  $0.94 * 1/3$  for the *ambiguous-uninformative* condition (Fig. 2). As

371 summarized in Fig. 2, the measured performances matched the assumptions in the contextual  
372 gaze following task very well (comparison by two-tailed t-tests, n.s.). This result indicates that  
373 the probability of identifying an object as a target was determined by the information provided  
374 by gaze direction and the verbal instruction.

375

376 *ROI localization.* To localize the GFP we contrasted *gaze following* with *color matching* trials in  
377 the first experiment. At the group level ( $N = 19$ ) this contrast (*gaze following* > *color matching*)  
378 yielded a patch of significantly larger activity for *gaze following* close to the pSTS in both  
379 hemispheres. The contrast maxima (blue spheres in Fig. 3, left column) were located at  $x, y, z = -$   
380  $57, -61, -1$  in the left and at  $x, y, z = 48, -67, -1$  in the right hemisphere. These locations closely  
381 match those known from previous studies, visualized as green and cyan spheres for comparison  
382 (Marquardt et al., 2017; Materna et al., 2008). In addition to the GFP, the *gaze following* > *color*  
383 *matching* contrast was also significant in a few more regions, not consistently seen as activated  
384 in previous work using the same paradigm (see extended data Fig. 3-1 for a list of all activated  
385 regions).

386 We localized the right hemispheric GFP in nine individual subjects (mean distance to group  
387 coordinates = 6.6 mm,  $SD = 3.1$  mm) and the left GFP in six subjects (mean distance = 7.7 mm;  
388  $SD = 1.4$  mm) (white spheres in Fig. 3, left column).

389 An analogous procedure was applied to localize the hLIP, using the contrast *gaze following* vs.  
390 *baseline fixation*. The location of maximum activation at the group level was found to be at  $x, y,$   
391  $z = 21, -67, 50$  (right) and  $x, y, z = -21, -67, 53$  (left) (blue spheres in Fig. 3, middle column) in  
392 good accordance with previous work on saccade related activity in the parietal cortex (Serenio et  
393 al., 2001) (Fig. 3, middle). We identified the hLIP regions bilaterally in all 19 subjects

394 individually with a mean distance of 13.4 mm ( $SD = 3.9$  mm) to the standard coordinates in the  
395 right hemisphere and 11.93 mm ( $SD = 3.7$  mm) in the left hemisphere (white spheres in Fig. 3,  
396 middle column).

397 In order to determine if BOLD activity in regions not delineated by the localizer experiment was  
398 modulated by the conditions of the contextual gaze following task, we contrasted activity in each  
399 of the ambiguous conditions with the *unambiguous* condition. The contrast *ambiguous-*  
400 *uninformative > unambiguous* was significant for a region in the inferior prefrontal cortex (Fig.  
401 3, right) whose group level maxima were found in slightly different locations in the two  
402 hemispheres, namely at  $x, y, z = -39, 11, 29$  in the left and  $x, y, z = 48, 20, 23$  in the right  
403 hemisphere (blue spheres), corresponding to the most lateral part of left BA 8 and upper right  
404 BA 44. In 15 subjects, we could delineate individual contrast locations (white spheres *ibid.*,  $SD$   
405 of individual locations (in mm): right  $x, y, z = 5, 6, 6$ ; left  $x, y, z = 5, 8, 6$ ). These individual  
406 locations were scattered around BA 44, BA 8 and BA 9 and henceforth we will refer to this  
407 region as the inferior frontal junction (IFJ). Contrasting *ambiguous-informative > unambiguous*  
408 yielded an overall weaker activation with only the right IFJ surviving a threshold of  $p < 0.001$ .

409 Weaker, albeit still significant *ambiguous-uninformative > unambiguous* contrasts were also  
410 found in the medial part of left BA 8 at  $x, y, z = -3, 11, 50$ , bilaterally in BA 6 at  $x, y, z = -21, -4,$   
411  $50$  and  $x, y, z = 24, -1, 50$  and at  $x, y, z = 36, 8, 47$  (right hemisphere) not far from the IFJ (cf.  
412 Extended data Fig. 3-1). Upon reversing the contrast, i.e. *unambiguous > ambiguous-*  
413 *informative/-uninformative*, we observed bihemispheric significance within BA 13 (insula), BA  
414 40, within the cingulate cortex (BA 24 and 31) and within BA 7 (all  $p < 0.001$ , and a minimum  
415 of 6 adjacent voxel, cf. Extended data Fig. 3-1).

416

417 *Time course of BOLD signals.* We wanted to know how the BOLD signal evolved over time,  
418 relative to the events of the trials in both of the a priori defined ROIs and the ad hoc identified  
419 IFJ. Fig. 4 shows the averaged time courses of the BOLD signal for each condition of the  
420 *contextual gaze following* experiment separately for the GFP and hLIP. This allows to present  
421 the response properties in the three conditions relative to baseline as well as their temporal  
422 pattern in relation to the *cue* and the *go*-signal not visible in the contrast-maps which only reflect  
423 the average signal of trials.

424 We performed two types of analyses to investigate the effects of context condition  
425 (*unambiguous*, *ambiguous-informative* and *ambiguous-uninformative*) on the BOLD activity; (1)  
426 permutation tests on each time point of the extracted BOLD signals (FDR corrected), (2)  
427 estimation of hierarchical models to infer BCIs of the time courses (cf. Extended data Fig. 6).

428 In the GFP, we observed two peaks throughout the trial, one at 10 s and the other one after 16.5  
429 s. Considering the latency of the BOLD signal of about 5 s we assume that the first peak is  
430 related to the onset of the *cue* (at 5 s) and the second to the *go*-signal at 10 s. For the GFP, we did  
431 not observe significant difference between any conditions at any time point.

432 The hLIP region depicted a similar two-peak pattern in response to the *cue* and the *go*-signal.  
433 Statistical analysis indicated that the BOLD response was significantly different between  
434 *unambiguous* and *ambiguous-uninformative* trials after 15 s (permutation test,  $q < 0.05$ , Fig. 4,  
435 bottom row) or after 13 s (BCI analysis, Fig. 6). The results of both analyses are in good  
436 agreement and – within the limits of the temporal resolution of the BOLD signal – suggest that  
437 the relevant event which caused the differentiation of the BOLD signal is the *go*-signal 10 s after  
438 trial onset rather than the *cue* 5 s after trial onset. Mind, that here, the model-based analysis  
439 method allowed a rejection of the null-hypothesis of no difference only for the left hemisphere

440 yet not for the right one. Since the pattern of the right hemisphere closely resembles the one of  
441 the left hemisphere and BCIs are only barely overlapping, we tend, however, to attribute this  
442 outcome to the low signal-to-noise ratio. This view is also supported by the decoding analysis  
443 (see below). There was no significant difference between the *ambiguous-informative* and  
444 *ambiguous-uninformative* conditions ( $q > 0.05$ ).

445 To rule out the possibility that the difference between the *unambiguous* and the *ambiguous-*  
446 *uninformative* condition was due to a larger number of saccades caused by the higher  
447 uncertainty, we performed a t-test on the number of saccades across subjects, which yielded no  
448 significant difference ( $p > 0.05$ ). To summarize, hLIP exhibited a significantly stronger activity  
449 in *ambiguous-uninformative* trials compared to *unambiguous* trials (at least in the left  
450 hemisphere) while this was not the case for the GFP.

451 Finally, we analyzed the time courses of the ROIs, identified in the exploratory whole brain  
452 analysis. Of these ROIs (cf. Fig. 3-1), only the IFJ survived the permutation test (Fig. 5 and 6).  
453 Notice that this analysis is not intended to compare the temporal average of the BOLD signal  
454 among the experimental conditions (which would be partly redundant to the activation map), but  
455 to get an idea about the temporal response characteristics during trials of the three conditions  
456 relative to baseline. Compared to GFP and hLIP, the condition-dependent BOLD signal in the  
457 IFJ exhibited a qualitatively different property: While the signal from the other two ROIs was  
458 modulated by the task during trials of all experimental conditions, the IFJ-signal did not exceed  
459 baseline signal during the *unambiguous* condition and was only modulated in the conditions  
460 comprising ambiguity. During the latter the signal was sustained at a higher level until the end of  
461 the trial. The permutation test yielded significant differences between the *unambiguous* and the  
462 *ambiguous-uninformative* conditions between 12 s and 17 s (left) and 12 s and 15 s (right)

463 (permutation test,  $q < 0.05$ , Fig. 5) or between 9.5 s and 16 s (BCI analysis, Fig. 6). Even if  
464 considering the 12 s-estimate as the onset of the differential, condition dependent modulation of  
465 the BOLD signal, the onset is too early to be related to the *go*-signal. Its association with the  
466 preceding event (the gaze and verbal cue 5 s earlier) seems more plausible. The profiles for the  
467 *ambiguous-informative* and *ambiguous-uninformative* conditions were statistically not different  
468 from each other. Summarizing, this analysis yielded that the BOLD signal from the IFJ  
469 significantly differentiated between conditions 3-4 s earlier than the hLIP BOLD signal,  
470 independent of the analysis method. With a temporal difference of 5 s between the main events  
471 of the task we tend to attribute the modulation of the IFJ-ROI to the *gaze/auditory cue* and the  
472 modulation of the hLIP-ROI the *go*-signal.

473

474 *Decoding of unambiguous vs. ambiguous-uninformative.* We also performed a decoding analysis  
475 by training a classifier on *ambiguous-uninformative* and *ambiguous* trials. The analysis yielded  
476 results which further support that only the parietal and the frontal regions but not the GFP are  
477 modulated by the contextual condition. The group-level t-map of classification accuracies (Fig.  
478 7) evince hotspots in accordance with previously described locations. One qualification is that  
479 the hotspots in the parietal cortex are slightly lateral to the original hLIP-ROIs whose locations  
480 were estimated based on the study by Sereno et al., 2001. An additional area in BA 19 that was  
481 not visible in the contrast map provided a decodable signal as well. In accordance with the time  
482 course analysis, no temporal area contributed to the decoding of conditions.

483 As a final step we extracted the mean accuracy values for voxels constituting the ROIs that were  
484 used for the time course analysis from each participant's accuracy map and compared the  
485 obtained distributions against a classifier performing at chance level. Here, the only distributions

486 that were significantly different from chance-level were those corresponding to the left and right  
487 IFJ-ROIs (Wilcoxon signed-rank test,  $p < 0.001$  (both hemispheres)). Neither the distributions  
488 stemming from the GFP-ROIs nor those from the original hLIP-ROIs were significantly different  
489 from chance level. Since we had found significant differences between *unambiguous and*  
490 *ambiguous-uninformative* conditions in the time-course analysis of hLIP this is surprising at the  
491 first glance. We think that the explanation lies in the procedure we used to determine the  
492 individual hLIP locations. This procedure involved using coordinates of a different study as seed  
493 coordinates and thus may have biased the localization of the hLIP (cf. Methods section). Indeed,  
494 the contrast shown in Fig. 3 (middle column) indicates that the activity within the parietal cortex  
495 is widespread and spans the superior parietal lobule and the adjoining sulci. Consequentially,  
496 postulating a confined hLIP in the context of the beta-contrasts may not be fully justified.  
497 However, the hotspot revealed by the decoding analysis may provide a more reliable estimate of  
498 the parietal area that is of relevance in our task and which is shifted in fronto-lateral directions as  
499 compared to the locations stemming from the beta-contrasts.

500 Irrespective of this qualification regarding the anatomical delineation of hLIP, both the time  
501 course analysis and the decoding analysis suggest differentially activated regions in the parietal  
502 and frontal cortex while the GFP does not come into play. Thus, we see our main finding to be  
503 confirmed by two independent analyses.

504

505

506 **Discussion**

507

508 In this study, we aimed at delineating the cortical regions that allow humans to single out objects  
509 being jointly attended to in case more than one object is hit by the other's gaze vector. To this  
510 end we ran 2 experiments. In experiment 1 we identified two distinct cortical regions (GFP and  
511 hLIP) known from previous works to be involved in processing the other's gaze direction and in  
512 shifting spatial attention respectively (Marquardt et al., 2017; Materna et al., 2008; Sereno et al.,  
513 2001). In experiment 2, subjects performed a *contextual gaze following* task in which they  
514 integrated gaze direction and auditory information to identify the objects attracting an avatar's  
515 attention. While BOLD activity of the GFP was not modulated by the informativeness of the  
516 auditory information, hLIP showed increased activity when the provided information was  
517 insufficient to specify the target. The BOLD contrast between the condition unambiguously  
518 specifying the targets and the two ambiguous conditions identified yet another area only  
519 involved in contextual gaze following, missed by the localizer paradigm due to the lack of the  
520 need to disambiguate object choices. This area exhibited a continuously elevated response if and  
521 only if the evidence about the target was low. Unlike the other two areas, IFJ did not show a  
522 general response to events of the trials in all experimental conditions; apart from an initial bump  
523 resembling the early part of the activity profiles during the two ambiguous conditions, its activity  
524 during *unambiguous* trials was close to or undistinguishable from baseline activity. This pattern  
525 of the IFJ activation suggests that the area is involved in the process of selecting the target if  
526 sensory stimuli would not unequivocally single out the target, i.e. IFJ may provide a top-down  
527 attention/ selection signal.

528 This study confirms previous findings that the GFP, which is located close to the pSTS, plays a  
529 major role in processing information on the others' gaze. Moreover, the present work shows that  
530 no matter if one or more potential target objects are hit by the other's gaze vector, the BOLD  
531 activity in the GFP remains the same. The need to differentiate between objects in case more  
532 than one lies on the gaze vector requires contributions from additional areas that exhibit  
533 differential activity. One of these areas, the hLIP in the posterior parietal lobe is also activated in  
534 the traditional, restricted gaze following paradigms in which the gaze hits one object only. It is  
535 established that the hLIP is necessary for the control of spatial attention (Corbetta and Shulman,  
536 2002). A role of hLIP is also supported by an additional decoding analysis that delineated voxels  
537 based on which the *unambiguous* and the *ambiguous-uninformative* conditions could be  
538 differentiated. This searchlight analysis revealed four hotspots, two of them in the posterior  
539 parietal cortex of both hemispheres, one in the left inferior prefrontal cortex and one in BA 19 of  
540 both hemispheres. The hotspots in the posterior parietal cortex turned out to be located slightly  
541 more rostral and inferior than the a-priori defined hLIP regions as identified based on locations  
542 taken from Sereno et al., (2001). We do not think that this spatial incongruence is surprising  
543 since the two localization approaches were based on different tasks (hLIP from experiment 1 and  
544 the hotspots from experiment 2) and, moreover, relied on different analysis methods (beta-  
545 contrasts versus decoding analysis). Hence, these two cortical spots, although overlapping or at  
546 least closely adjacent may not necessarily play the same role in gaze following. In any case, our  
547 statistical analyses showed that both, hLIP and the hotspots were sensitive to the experimental  
548 conditions. Future work may investigate the exact roles of both regions regarding contextual  
549 gaze following. The searchlight analysis also revealed a significant contribution of an area in BA  
550 19, not identified by the original search for significant BOLD contrasts. BOLD activity in

551 extrastriate BA 19 is typically seen in tasks that involve varying contributions to complex visual  
552 processing. Its functional role has remained elusive. In view of the fact that BA 19 of nonhuman  
553 primates maintains bidirectional connections with BA 7 (Rockland and Pandya, 1979) one may  
554 speculate that a significant contribution to the decoding of conditions decoding may arise from  
555 feedback connections originating in hLIP.

556 Work on monkey area LIP, arguably homologous to hLIP, has suggested that this area  
557 constitutes a priority or saliency map providing a representation of the environment that  
558 highlights locations that serve as attractors of attention. The saliency map may be modulated by  
559 bottom-up sensory cues, symbolic cues or gaze cues (Bisley and Goldberg, 2010; Walther and  
560 Koch, 2006). The latter is suggested by single unit recordings from area LIP. Many LIP neurons  
561 respond to the appearance of a gaze cue provided the gazed at location lies within the neuron's  
562 receptive field (Shepherd et al., 2009). Spatial selectivity for gazed at locations and objects at  
563 these locations is also exhibited by many neurons in the monkey GFP (Ramezanzpour and Thier,  
564 2020). This selectivity suggests that the priority map in LIP might draw on input from the GFP.  
565 The yoked activation of hLIP/LIP and the GFP in BOLD imaging studies of gaze following is in  
566 principle in accordance with this scenario (Marquardt et al., 2017; Materna et al., 2008; Shepherd  
567 et al., 2009). However, the poor temporal resolution of the BOLD signals does not allow us to  
568 critically test if the assumed direction of information flow holds true. In any case, bidirectional  
569 projections are known to connect monkey area LIP and parts of the STS (Seltzer and Pandya,  
570 1994). One well-established pathway links area LIP and PITd, an area in the lower STS,  
571 probably close to the GFP, known to contribute to the maintenance of sustained attention (Sani et  
572 al., 2019; Stemmann and Freiwald, 2016). Yet, the anatomical data available does not allow us to  
573 decide if the GFP does indeed contribute to this fiber bundle.

574 In the present study the BOLD signal evoked by gaze following in hLIP was overall much  
575 stronger than in the GFP. Moreover, unlike the GFP signal, it exhibited a dependence on the  
576 conditions of the *contextual gaze following* experiment. Higher activity was associated with the  
577 *ambiguous-informative* and the *ambiguous-uninformative* conditions, both associated with  
578 unresolved uncertainty about the object requiring a decision of the participant that could only  
579 partially be based on information provided by the cue. Why should a region thought to  
580 coordinate spatial shifts of attention show an influence of target ambiguity, i.e. the need to  
581 choose between several potential targets? One possible answer may be that the higher hLIP  
582 activity reflects an increased attentional load. More specifically, increased uncertainty in  
583 ambiguous trials may have prompted more shifts of attention from one object to the other in an  
584 attempt to resolve the ambiguity. Although we found no difference in the number of exploratory  
585 saccades after the *go*-signal across conditions, we cannot rule out that participants covertly  
586 shifted attention between targets in ambiguous trials more than in the other trials. However, a  
587 more parsimonious explanation could be that hLIP constitutes a neural substrate for making  
588 decisions under uncertainty independent of the attentional load as suggested by several studies  
589 such as (Vickery and Jiang, 2009).

590 The BOLD signal in the area we identified as the IFJ (between premotor cortex (BA 6), BA 44  
591 and BA 8) exhibited a dependency on condition as well. This result is suggested by the BOLD  
592 contrast (Fig. 3, right column), the time course analysis (Fig. 5 and 6) as well as by the decoding  
593 analysis (Fig. 7 and 7-1). However, the time course analysis revealed a fundamental difference  
594 compared to response profiles of BOLD activity in hLIP or the GFP. Sustained activity could  
595 only be observed in trials of the two *ambiguous* conditions, i.e. when the participants needed to  
596 make decisions under sensory uncertainty. This suggests that the condition dependency of the IFJ

597 signal may be a consequence of shifts of attention between the two object categories, houses and  
598 hands. This interpretation draws on a MEG-fMRI study that demanded the allocation of attention  
599 to distinct classes of visual objects such as faces and spatial scenes (Baldauf and Desimone,  
600 2014). Depending on the object of attention, gamma band activity in the IFJ was synchronized  
601 either with the fusiform face area (FFA) or the parahippocampal place area (PPA). Additional  
602 support for this view comes from spatial cueing paradigms, which suggest that the IFJ primarily  
603 supports transient attentional processes, such as covert attentional shifts (Asplund et al., 2010;  
604 Tamber-Rosenau et al., 2018). We speculate that the time course of activity in the IFJ reflects the  
605 coordination of covert shifts of attention until the choice for the saccade target is made. In  
606 unambiguous trials, the lack of ambiguity allows fast decisions and since no attentional shifts are  
607 necessary the IFJ is not required.

608 The functional characteristics of the GFP, hLIP and the IFJ attribute complementary functions to  
609 each area which, in sum, allows gaze following under sensory ambiguity. We propose that  
610 information on the direction of the other's gaze is provided by the GFP and modulates the  
611 saliency map generated by area hLIP such that spatial positions in the direction of the gaze  
612 vector are highlighted. In this situation the choice of which of the possible objects is the most  
613 relevant one requires the resolution of uncertainty which is accomplished by the IFJ. In this  
614 scenario the intersection between the spatial information provided by the GFP-hLIP complex and  
615 the object-based information provided by the IFJ singles out one object that will then become the  
616 target of the observer's gaze following response, elicited by the hLIP.

617 Several points need to be addressed by future work in order to test and to further refine this  
618 concept. As a first step, it will be necessary to investigate the temporal interplay between these  
619 regions in an attempt to establish causal interactions in order to critically test the model. Our

620 hypothesis assumes that the IFJ has a leading role in processing information on competing  
621 objects on the gaze vector, resolving the uncertainty as to which one the target is. The conclusion  
622 that IFJ has a leading role in the disambiguation of the object set is primarily based on the fact  
623 that ambiguity related information arises first in IFJ and only later in hLIP. Yet, we cannot rule  
624 out that this sequence might be an artifact of region-specific differences in the statistical power  
625 of the BOLD time course analysis, eventually in conjunction with region specific differences in  
626 the variability of BOLD signal latencies.

627 To summarize, our results suggest a fronto-temporo-parietal network for geometric gaze  
628 following and the allocation of joint attention. While the GFP seems to have a leading role in  
629 selecting objects identified by the other's gaze vector, the IFJ seems to play a central role in  
630 disambiguating object choices in case more than one object may be hit by the vector. Finally, the  
631 hLIP acts as a priority map, highlighting the spatial location of the target object based on the  
632 focus of attention, and contribute to the execution of gaze shifts.

633

634

635 **References**

- 636 Asplund CL, Todd JJ, Snyder AP, Marois R (2010) A central role for the lateral prefrontal cortex  
637 in goal-directed and stimulus-driven attention. *Nature Neuroscience* 13:507–512.
- 638 Atabaki A, Marciniak K, Dicke PW, Thier P (2015) Assessing the precision of gaze following  
639 using a stereoscopic 3D virtual reality setting. *Vision Res* 112:68–82.
- 640 Baldauf D, Desimone R (2014) Neural Mechanisms of Object-Based Attention. *Science*  
641 344:424–427.
- 642 Baron-Cohen S (1995) *Mindblindness: An essay on autism and theory of mind*, *Mindblindness:*  
643 *An essay on autism and theory of mind*. Cambridge, MA, US: The MIT Press.
- 644 Baron-Cohen S (1994) How to build a baby that can read minds: Cognitive mechanisms in mind  
645 reading. *Curr Psychol Cogn* 13:513–552.
- 646 Benjamini Y, Hochberg Y (1995) Controlling the False Discovery Rate: A Practical and  
647 Powerful Approach to Multiple Testing. *Journal of the Royal Statistical Society Series B*  
648 (Methodological) 57:289–300.
- 649 Bisley JW, Goldberg ME (2010) Attention, Intention, and Priority in the Parietal Lobe. *Annual*  
650 *Review of Neuroscience* 33:1–21.
- 651 Bock SW, Dicke PW, Thier P (2008) How precise is gaze following in humans? *Vision Res*  
652 48:946–957.
- 653 Butterworth G, Jarrett N (1991) What minds have in common is space: Spatial mechanisms  
654 serving joint visual attention in infancy. *British journal of developmental psychology*  
655 9:55–72.
- 656 Corbetta M, Shulman GL (2002) Control of goal-directed and stimulus-driven attention in the  
657 brain. *Nature Reviews Neuroscience* 3:201–215.
- 658 Emery NJ (2000) The eyes have it: The neuroethology, function and evolution of social gaze.  
659 *Neuroscience & Biobehavioral Reviews* 24:581–604.
- 660 Hebart MN, Georgen K, Haynes J-D (2015) The Decoding Toolbox (TDT): A versatile software  
661 package for multivariate analysis of functional imaging data. *Frontiers in*  
662 *Neuroinformatics* 8.
- 663 Hollander G de, Knapen T (2017) nideconv. Vrije Universiteit and the Spinoza Centre for  
664 Neuroimaging.
- 665 Kobayashi H, Kohshima S (1997) Unique morphology of the human eye. *Nature* 387:767–768.
- 666 Lacadie CM, Fulbright RK, Rajeevan N, Constable RT, Papademetris X (2008) More accurate  
667 Talairach coordinates for neuroimaging using non-linear registration. *NeuroImage*  
668 42:717–725.
- 669 Laube I, Kamphuis S, Dicke PW, Thier P (2011) Cortical processing of head- and eye-gaze cues  
670 guiding joint social attention. *Neuroimage* 54:1643–1653.
- 671 Liuzza MT, Cazzato V, Vecchione M, Crostella F, Caprara GV, Aglioti SM (2011) Follow My  
672 Eyes: The Gaze of Politicians Reflexively Captures the Gaze of Ingroup Voters. *PLOS*  
673 *ONE* 6:e25117.
- 674 Marquardt K, Ramezanzpour H, Dicke PW, Thier P (2017) Following Eye Gaze Activates a Patch  
675 in the Posterior Temporal Cortex That Is Not Part of the Human “Face Patch” System.  
676 *eNeuro* 4:1–10.

- 677 Materna S, Dicke PW, Thier P (2008) Dissociable Roles of the Superior Temporal Sulcus and  
678 the Intraparietal Sulcus in Joint Attention: A Functional Magnetic Resonance Imaging  
679 Study. *J Cogn Neurosci* 20:108–119.
- 680 Ramezanzpour H, Thier P (2020) Decoding of the other’s focus of attention by a temporal cortex  
681 module. *PNAS*.
- 682 Rockland KS, Pandya DN (1979) Laminar origins and terminations of cortical connections of the  
683 occipital lobe in the rhesus monkey. *Brain Res* 179:3–20.
- 684 Sani I, McPherson BC, Stemmann H, Pestilli F, Freiwald WA (2019) Functionally defined white  
685 matter of the macaque monkey brain reveals a dorso-ventral attention network. *eLife*  
686 8:e40520.
- 687 Seltzer B, Pandya DN (1994) Parietal, temporal, and occipita projections to cortex of the  
688 superior temporal sulcus in the rhesus monkey: A retrograde tracer study. *J Comp Neurol*  
689 343:445–463.
- 690 Sereno MI, Pitzalis S, Martinez A (2001) Mapping of Contralateral Space in Retinotopic  
691 Coordinates by a Parietal Cortical Area in Humans. *Science* 294:1350–1354.
- 692 Shepherd SV, Klein JT, Deaner RO, Platt ML, Shepard RN (2009) Mirroring of attention by  
693 neurons in macaque parietal cortex. *Proc Natl Acad Sci* 106:9489–9494.
- 694 Stan Development Team (2018) PyStan: the Python interface to Stan.
- 695 Stemmann H, Freiwald WA (2016) Attentive Motion Discrimination Recruits an Area in  
696 Inferotemporal Cortex. *J Neurosci* 36:11918–11928.
- 697 Tamber-Rosenau BJ, Asplund CL, Marois R (2018) Functional dissociation of the inferior  
698 frontal junction from the dorsal attention network in top-down attentional control. *Journal*  
699 *of Neurophysiology* 120:2498–2512.
- 700 Teufel C, Alexis DM, Clayton NS, Davis G (2010) Mental-state attribution drives rapid,  
701 reflexive gaze following. *Atten Percept Psychophys* 72:695–705.
- 702 Teufel C, Alexis DM, Todd H, Lawrance-Owen AJ, Clayton NS, Davis G (2009) Social  
703 Cognition Modulates the Sensory Coding of Observed Gaze Direction. *Curr Biol*  
704 19:1274–1277.
- 705 Vickery TJ, Jiang YV (2009) Inferior Parietal Lobule Supports Decision Making under  
706 Uncertainty in Humans. *Cereb Cortex* 19:916–925.
- 707 Walther D, Koch C (2006) Modeling attention to salient proto-objects. *Neural Networks, Brain*  
708 *and Attention* 19:1395–1407.
- 709  
710

711  
712 **Fig. 1. Contextual gaze following task.** An avatar appeared in the center of the screen together with four linearly  
713 arranged sets of objects (houses and hands). After a baseline fixation period, the portrait's gaze shifted towards one  
714 specific target object simultaneously with an auditory contextual instruction specifying the object class of the target  
715 (hand or house) or not, i.e. remaining uninformative ("none"). While maintaining fixation, subjects needed to decide  
716 on the target and make a saccade to the chosen target after a *go*-signal indicated by the disappearance of the fixation  
717 dot.

718

719 **Fig. 2. Behavioral performance.** Left: Boxplots (black and gray) showing the percentage of correct response in the  
720 localizer paradigm (dashed line depicts chance level performance, see Fig. 2-1 for a description of the localizer  
721 paradigm). Right: Plots of correct responses in the contextual gaze following paradigm (weighted mean performance  
722 and weighted *SD*, dashed lines depict expected performance; blue: *unambiguous*, yellow: *ambiguous-informative*,  
723 salmon: *ambiguous-uninformative*).

724

725 **Fig. 3. Activation maps emphasizing the ROIs.** Left column: contrast *gaze following* > *color matching* (localizer  
726 paradigm) used to identify the GFP. Blue dots mark maximum activation on the group level closest to locations  
727 taken from literature (green (Marquardt et al., 2017) and cyan (Materna et al., 2008) dots), white dots mark the  
728 maximum activation of those locations which were identifiable on the individual level. Middle column: contrast  
729 *gaze following* > *baseline fixation* (localizer paradigm) used to identify saccade-related activity in the hLIP closest  
730 to location taken from (Serenio et al., 2001) (cyan dot). Blue and white dots mark again, group level and individual  
731 coordinates; Right column: *ambiguous-uninformative* > *unambiguous* (contextual gaze following paradigm). Blue  
732 and white dots mark the group level and individual locations of the maximum IFJ-activity. See extended data Fig. 3-  
733 1 for tabular form of all activated regions of each contrast and Fig. 3-2/3/4 for the respective contrast maps.

734

735 **Fig. 4. Time courses of activation in the GFP and the hLIP.** Time course of mean percent signal change in the  
736 contextual gaze following experiment in areas identified in the localizer experiment (error bars are SEM). Areas in  
737 which conditions showed significant differences are shaded (permutations test,  $q < 0.05$ ). See Fig. 6 for the mode-  
738 based analysis.

739

740 **Fig. 5. Time courses of activation in the IFJ.** Time course of mean percent signal change during the contextual  
741 gaze following experiment of the IFJ (error bars are SEM). Areas in which conditions showed significant differences  
742 are shaded (permutations test,  $q < 0.05$ ). See Fig. 6 for the mode-based analysis.

743 **Fig. 6:** Bayesian hierarchical model of time courses. Shaded areas comprise 95% credible intervals.

744 **Fig. 7:** Group level result of the searchlight decoding analysis (t-map of classification accuracies,  $p < 0.001$ ). Fig. 7-  
745 1 shows the distributions of individual accuracies for the ROIs used in the time course analysis.

746 **Fig. 3-1:** List of activated brain areas.

747 **Fig. 3-2:** Contrast map *gaze following* > *color matching* ( $p < 0.001$ , cluster size > 6 voxels).

748 **Fig. 3-3:** Contrast map *gaze following* > *baseline* ( $p < 0.001$ , cluster size > 6 voxels).

749 **Fig. 3-4:** Contrast map *ambiguous-uninformative* > *unambiguous* ( $p < 0.001$ , cluster size > 6 voxels).

750 **Fig. 7-1:** Classification accuracy distributions across participants for the ROIs used in the time course analysis.  
751 Asterisks marking distributions significantly different from a classifier performing at chance level (Wilcoxon  
752 signed-rank test,  $p < 0.001$ ).

753

754

755

756

Baseline Fixation



5 sec

Cue



5 sec

Go



2 sec

Time











