

1 **Title:** Temporo-Parietal Cortex Involved in Modeling One’s Own and Others’ Attention

2 **Authors:** Arvid Guterstam\*, Branden J Bio, Andrew I Wilterson, Michael SA Graziano

3 **Affiliation:** Department of Psychology, Princeton University, Princeton, NJ 08544

4 **\*Corresponding author:** arvidg@princeton.edu

5

6 **Abstract:** In a traditional view, in social cognition, attention is equated with gaze and  
7 people track attention by tracking other people’s gaze. Here we used fMRI to test  
8 whether the brain represents attention in a richer manner. People read stories describing  
9 an agent (either oneself or someone else) directing attention to an object in one of two  
10 ways: either internally directed (endogenous) or externally induced (exogenous). We  
11 used multivoxel pattern analysis to examine how brain areas within the theory-of-mind  
12 network encoded attention type and agent type. Brain activity patterns in the left  
13 temporo-parietal junction (TPJ) showed significant decoding of information about  
14 endogenous versus exogenous attention. The left TPJ, left superior temporal sulcus  
15 (STS), precuneus, and medial prefrontal cortex (MPFC) significantly decoded agent type  
16 (self versus other). These findings show that the brain constructs a rich model of one’s  
17 own and others’ attentional state, possibly aiding theory of mind.

18

19 **Impact statement:** This study used fMRI to show that the human brain encodes other  
20 people’s attention in enough richness to distinguish whether that attention was directed  
21 exogenously (stimulus-driven) or endogenously (internally driven).

## 1 **Introduction**

2       Reconstructing someone else's attentional state is of central importance in theory of  
3 mind (Baron-Cohen, 1997; Calder et al., 2002; Graziano, 2013). By identifying the object  
4 of someone else's attention, and having some intuitive understanding of the complex  
5 dynamics and consequences of attention, one can reconstruct at least some of the other  
6 person's likely thoughts, intentions, and emotions, and make predictions about that  
7 person's behavior. Almost all work on how people reconstruct the attention of others has  
8 focused on gaze direction. For example, the human eye has a high contrast between pupil  
9 and sclera, possibly an adaptation for better gaze tracking (Kobayashi and Kohshima,  
10 1997). The superior temporal sulcus in monkeys and humans may contain specialized  
11 neural circuitry for processing gaze direction (Hoffman and Haxby, 2000; Perrett et al.,  
12 1985; Puce et al., 1998; Wicker et al., 1998). Seeing a face gaze at an object  
13 automatically draws one's own attention to the object (Friesen and Kingstone, 1998;  
14 Frischen et al., 2007). These and other findings show the importance of reconstructing  
15 gaze direction in social cognition.

16       To be adaptive in aiding theory of mind, however, a model of attention should be  
17 far more than a vector indicating gaze direction. We previously suggested that the human  
18 brain constructs a rich, dynamic, and predictive model of other people's attention  
19 (Graziano, 2019, 2013; Graziano and Kastner, 2011). The model should contain  
20 information about different types of attention, about the rapidity or sluggishness with  
21 which attention tends to move from item to item, about how external factors such as  
22 salience and clutter are likely to affect a person's attention, and about how attention  
23 profoundly affects thought, memory, and behavior. In the proposal, that deeper model is

1 constrained by incoming information, including gaze direction. However, other cues can  
2 also constrain the model. People rely on the other person's body posture, on cues in the  
3 surrounding environment, on speech, and on social context. For example, blind people  
4 must be able to build models of other people's attention without seeing the other person's  
5 eyes. Likewise, during a phone conversation, we cannot see the other person and yet we  
6 intuitively understand whether that person is attending to what we have said or is  
7 distracted by her own words or by a salient event on her end of the line.

8         Several recent experiments provide evidence for an automatically constructed  
9 model of the attention of others that may go beyond merely registering gaze direction  
10 (Guterstam et al., 2019; Guterstam and Graziano, 2020; Kelly et al., 2014; Pesquita et al.,  
11 2016; Vernet et al., 2019). For example, Pesquita et al. (2016) found that when  
12 participants watch an actor in a video attending to an object, the participants implicitly  
13 distinguish between whether the actor's attention was drawn to the object exogenously  
14 (bottom-up, or stimulus-driven attention), or whether the actor endogenously shifted  
15 attention to the object (top-down, or internally driven attention). The ability to distinguish  
16 between someone else's exogenous and endogenous attention is one example of how  
17 people may construct a rich, dynamic model of other people's attention beyond merely  
18 encoding gaze direction or identifying the object of attention.

19         Inspired by the vignette-style tasks widely used in studies on theory of mind  
20 (Fletcher et al., 1995; Gallagher et al., 2000; Happé, 1994; Saxe and Kanwisher, 2003;  
21 Vogeley et al., 2001), in the present study, we used functional magnetic resonance  
22 imaging (fMRI) and multi-voxel pattern analysis (MVPA) to study brain activity in  
23 participants while they read brief stories about people's attention. Some of the stories

1 implied that attention was being attracted exogenously (“Kevin walks into his closet and  
2 notices the bright red tie...”) and some stories implied that attention was being directed  
3 endogenously (“Kevin walks into his closet and looks for the bright red tie...”). We also  
4 included analogous stories written in the first person, casting the subject of the  
5 experiment as the agent (“You walk into your closet and notice the bright red tie...”). The  
6 study therefore used a 2X2 design (exogenous versus endogenous attention X self agent  
7 versus other agent). Finally, we included a fifth, control condition, consisting of  
8 nonsocial stories in which the agent was replaced by an inanimate object, that, like  
9 attention, has a source and a target, such as a camera or a light source (“In a closet, a light  
10 shines on a red tie...”).

11 We made four predictions. Our first, central prediction was inspired by the Pesquita  
12 et al. (2016) study described above. We hypothesized that participants would encode the  
13 type of attention in the story (exogenous versus endogenous), and that this encoding  
14 would be evident in some subset of the areas classically involved in theory of mind.  
15 Previous experiments on theory of mind typically recruited a network of cortical areas  
16 including the temporoparietal junction (TPJ), the superior temporal sulcus (STS), the  
17 medial prefrontal cortex (MPFC), and the precuneus (Gallagher et al., 2000; Saxe and  
18 Kanwisher, 2003; van Veluw and Chance, 2014; Vogeley et al., 2001). We therefore  
19 predicted that the exogenous-versus-endogenous distinction would be significantly  
20 encoded in some subset of these areas. In particular, we anticipated that the TPJ might  
21 show the clearest evidence of encoding information about the type of attention, since  
22 previous experiments pointed to the TPJ as contributing to encoding other people’s  
23 attentional state when participants looked at faces gazing toward objects (Igelström et al.,

1 2016; Kelly et al., 2014). This first prediction, that the social cognition network will  
2 encode the exogenous-versus-endogenous distinction, represents the main, novel  
3 contribution of this study.

4       Second, we predicted that participants would encode information about the agent in  
5 the story (self versus other), and that this encoding would again be evident in some subset  
6 of the areas classically involved in theory of mind. Self-versus-other encoding has been  
7 examined in previous studies, and found to be reflected in the theory-of-mind network  
8 (e.g., Northoff et al., 2006; Ochsner et al., 2004; Passingham et al., 2010; Qin and  
9 Northoff, 2011; van Veluw and Chance, 2014). This second prediction represents a test of  
10 whether our present paradigm, using subtle wording differences between similar  
11 sentences, can produce results consistent with previous findings.

12       Third, we predicted that participants would encode information associated with the  
13 interaction between the two factors. We predicted that at least some subset of the areas in  
14 the theory-of-mind network may encode the type of attention (exogenous versus  
15 endogenous) to a different extent in self-related stories as compared to other-related  
16 stories.

17       Fourth and finally, we tested for brain regions that encoded the distinction between  
18 social stories (with human agents) and nonsocial stories (with only non-agent objects).  
19 We predicted that this social-versus-nonsocial encoding would again be evident in the  
20 same network of brain regions noted above, that are known to be involved in theory of  
21 mind. This final analysis served as a control to check on the validity of the story stimuli  
22 and confirm that they engaged social cognition as expected.

23

## 1 **Methods**

### 2 *Subjects*

3           Thirty-two healthy human volunteers (12 females, 30 righthanded, aged 18-52,  
4 normal or corrected to normal vision) participated in the study, based on the sample size  
5 used in a previous, related study (Guterstam et al., 2020). Subjects were recruited either  
6 from a paid subject pool, receiving 40 USD for participation, or from among Princeton  
7 undergraduate students, who received course credits as compensation. In the subject  
8 recruitment material, the experiment was described as a “Reading Comprehension  
9 Study.” All subjects provided informed consent and all procedures were approved by the  
10 Princeton Institutional Review Board.

11

### 12 *Experimental setup*

13           Before scanning, subjects were instructed and then shown three sample trials  
14 (which were not part of the stories presented in the subsequent experiment) on a laptop  
15 computer screen. All subjects gave the correct response to all three trials on the first try,  
16 indicating they had understood the instructions adequately. During scanning, the subjects  
17 laid comfortably in a supine position on the MRI bed. Through an angled mirror mounted  
18 on top of the head coil, they viewed a translucent screen approximately 80 cm from the  
19 eyes, on which visual stimuli were projected with a Hyperion MRI Digital Projection  
20 System (Psychology Software Tools, Sharpsburg, PA, USA) with a resolution of 1920 x  
21 1080 pixels. A PC running MATLAB (MathWorks, Natick, MA, USA) and the  
22 Psychophysics Toolbox (Brainard, 1997) was used to present visual stimuli. A right hand  
23 5-button response unit (Psychology Software Tools Celeritas, Sharpsburg, PA, USA) was

1 strapped to the subjects' right wrist. Subjects used the right index finger button to  
2 indicate a true response, and the right middle finger to indicate a false response during the  
3 probe phase of each trial.

4

#### 5 *Experimental conditions and stimuli*

6 Five experimental conditions were included. Subjects were presented with short  
7 stories (2-3 sentences, average word count = 24) describing a scene in which an agent,  
8 which was either the subject him-/herself (self) or another person (other), directed  
9 attention to something in the external world endogenously (e.g., "X is attentively looking  
10 for Y") or exogenously (e.g., "X's attention is captured by Y"). These four conditions  
11 made up a 2 x 2 factorial design: attention type (endogenous versus exogenous) X agent  
12 (self versus other). In addition, we included a baseline condition featuring stories in  
13 which the agent was substituted by a non-human object. In each trial, after a 9 – 11 s  
14 inter-trial interval, the story was presented for 10 s in easily readable, white text on a  
15 black background, at the center of the screen, after which a probe statement was shown  
16 for 4 s, to which the subjects responded either true or false by button press. See Figure 1  
17 for details, and SI Data S1 for all stories.

18 Each subject ran 100 trials and thus saw 100 stories: 80 social stories and 20 non-  
19 social control stories. The 80 social stories were constructed as follows. We began with  
20 80 unique short stories. For each story, four versions were constructed, one for each of  
21 the factorial conditions (Figure 1B). To keep the story versions as semantically similar as  
22 possible, we made minimal changes to the wordings. To distinguish the self and other  
23 versions, we substituted the word "you" with a name (e.g., "Karen") and the word "your"

1 with “his” or “her”. The names in the stories were selected from a list of the 100 most  
2 popular given names for male and female babies born during the years 1919-2018 in the  
3 United States, which is published by the Social Security Administration  
4 (<https://www.ssa.gov/oact/babynames/decades/century.html>). Half of the names were  
5 masculine, half feminine. To distinguish the endogenous and exogenous story versions,  
6 we used different wording for the part of the story where the agent (X) is related to the  
7 object (Y). In the endogenous versions, we used formulations such as: “X is trying to find  
8 Y,” “X is trying to spot Y,” or “X is looking attentively for Y.” In the exogenous  
9 versions, we used formulations such as: “X’s eyes are drawn to Y,” “X’s gaze is captured  
10 by Y,” or “X’s attention is captured by Y.” We matched the average number of words  
11 across all four conditions (24 words). The number of stories that included the words  
12 “attention” or “attentively” was balanced between the endogenous and exogenous  
13 categories (43 stories in each). Among the 80 stories, for each subject, 20 were randomly  
14 selected to be used in the endogenous-self version; 20 in the endogenous-other version;  
15 20 in the exogenous-self version; 20 in the exogenous-other version. Thus, for the  
16 example story shown in Figure 1B, each subject saw only one of the four versions. In this  
17 manner, each subject saw 80 social stories, 20 of each type, balanced for as many  
18 properties as possible other than the two factors that were manipulated.

19 Finally, we constructed 20 additional stories for the non-social control condition  
20 (Figure 1B). To keep the control stories as semantically similar as possible to the social  
21 stories, we based them on a subset of the 80 original stories. Crucially, the agent in the  
22 original story was substituted with a non-human object, such as a camera or a spotlight,  
23 that has a source and a target just as attention does. For instance, the original story, “You



1 are in a bike shop, and numerous bikes hang on one of the walls. You are attentively  
2 looking for that red Italian sports bike,” was adapted to the non-social condition by  
3 substituting the agent with a spotlight: “In a bike shop, on one of the walls, hangs  
4 numerous bikes. A bright spotlight is shining on a red Italian sports bike.” The average  
5 number of words of the non-social stories (24 words, standard deviation = 3) was  
6 matched with the attention stories.

7         The purpose of the probe statement at the end of each trial was to ensure that  
8 subjects carefully read the stories. Each statement described one detail of the preceding  
9 story that could be either true or false. We restricted the probe statements to the spatial  
10 context of the story (place probe: e.g., “Emma is on a bus”) or the object being described  
11 (object probe: e.g., “The Van Gogh painting has sunflowers”) in order to avoid alerting  
12 subjects to the focus of the experiment on theory of mind and attention. Half of the probe  
13 statements were place probes and half object probes. Within both the place and the object  
14 probes, half were true and half were false. The probe was on screen for 4 s, during which  
15 subjects were required to indicate whether the statement was true or false by button press.

16         The experiment consisted of 10 runs of approximately 4 min each. In each run,  
17 the 5 conditions were repeated 2 times, yielding a total of 10 trials per run. The trial order  
18 was randomized, with the limitation that two consecutive trials could not belong to the  
19 same condition. Each run included 18 s of baseline before the onset of the first trial and  
20 12 s of baseline after the offset of the last trial.

21

## 22 *Post-scan questionnaire*

23         At the end of the scanning session, subjects were asked what they thought the

1 purpose of the experiment was and what they thought it was testing.

2

### 3 *fMRI data acquisition*

4 Functional imaging data were collected using a Siemens Prisma 3T scanner  
5 equipped with a 64-channel head coil. Gradient-echo T2\*-weighted echo-planar images  
6 (EPI) with blood-oxygen dependent (BOLD) contrast were used as an index of brain  
7 activity (Logothetis et al., 2001). Functional image volumes were composed of 54 near-  
8 axial slices with a thickness of 2.5 mm (with no interslice gap), which ensured that the  
9 entire brain excluding cerebellum was within the field-of-view in all subjects (54 x 78  
10 matrix, 2.5 mm x 2.5 mm in-plane resolution, TE = 30 ms, flip angle = 80°).  
11 Simultaneous multi-slice (SMS) imaging was used (SMS factor = 2). One complete  
12 volume was collected every 2 s (TR = 2000 ms). A total of 1300 functional volumes were  
13 collected for each participant, divided into 10 runs (130 volumes per run). The first three  
14 volumes of each run were discarded to account for non-steady-state magnetization. A  
15 high-resolution structural image was acquired for each participant at the end of the  
16 experiment (3D MPRAGE sequence, voxel size = 1 mm isotropic, FOV = 256 mm, 176  
17 slices, TR = 2300 ms, TE = 2.96 ms, TI = 1000 ms, flip angle = 9°, iPAT GRAPPA = 2).  
18 At the end of each scanning session, matching spin echo EPI pairs (anterior-to-posterior  
19 and posterior-to-anterior) were acquired for blip-up/blip-down field map correction.

20

### 21 *fMRI preprocessing*

22 Results included in this manuscript come from preprocessing performed using  
23 FM RIPREP version 1.2.3 (Esteban et al., 2019) (RRID:SCR\_016216), a Nipype

1 (Gorgolewski et al., 2011) (RRID:SCR\_002502) based tool. Each T1w (T1-weighted)  
2 volume was corrected for INU (intensity non-uniformity) using N4BiasFieldCorrection  
3 v2.1.0 (Tustison et al., 2010) and skull-stripped using antsBrainExtraction.sh v2.1.0  
4 (using the OASIS template). Spatial normalization to the ICBM 152 Nonlinear  
5 Asymmetrical template version 2009c (Fonov et al., 2009) (RRID:SCR\_008796) was  
6 performed through nonlinear registration with the antsRegistration tool of ANTs v2.1.0  
7 (Avants et al., 2008) (RRID:SCR\_004757), using brain-extracted versions of both T1w  
8 volume and template. Brain tissue segmentation of cerebrospinal fluid (CSF), white-  
9 matter (WM) and gray-matter (GM) was performed on the brain-extracted T1w using fast  
10 (Zhang et al., 2001) (FSL v5.0.9, RRID:SCR\_002823).

11 Functional data was slice time corrected using 3dTshift from AFNI v16.2.07  
12 (Cox, 1996) (RRID:SCR\_005927) and motion corrected using mcflirt (FSL v5.0.9)  
13 (Jenkinson et al., 2002). This was followed by co-registration to the corresponding T1w  
14 using boundary-based registration (Greve and Fischl, 2009) with six degrees of freedom,  
15 using flirt (FSL). Motion correcting transformations, BOLD-to-T1w transformation and  
16 T1w-to-template Montreal Neurological Institute (MNI) warp were concatenated and  
17 applied in a single step using antsApplyTransforms (ANTs v2.1.0) using Lanczos  
18 interpolation.

19 Many internal operations of FM RIPREP use Nilearn (Abraham et al., 2014)  
20 (RRID:SCR\_001362]) principally within the BOLD-processing workflow. For more  
21 details of the pipeline see <https://fmripred.readthedocs.io/en/latest/workflows.html>.

22

23 *Testing prediction 1*

1           The purpose of the first analysis was to determine whether the brain encoded  
2 information concerning the type of attention (endogenous or exogenous) present in the  
3 stories. For this analysis, we used MVPA, which tests whether patterns of brain activity  
4 can be used to decode the distinction between two conditions. It is a more sensitive  
5 analysis than the more common, simple subtraction methods. The reason for using this  
6 sensitive measure is that the difference between exogenous and endogenous trial types  
7 was extremely subtle. Both trial types engaged social cognition, and therefore might  
8 cancel each other out in a simple subtraction. The stimuli were nearly identical, differing  
9 only in a few words that indicated the type of attention used by the agent in the story. In  
10 addition, the type of attention featured in the story was irrelevant to the task performed by  
11 the subject. To accommodate the subtlety of the distinction between conditions, we  
12 designed the study to use MVPA. We hypothesized that with MVPA, brain activity  
13 would carry information about the endogenous versus exogenous distinction; and that  
14 decoding would be evident in regions of interest (ROIs) within the network of areas  
15 typically found to be involved in social cognition, especially within the TPJ.

16           We defined our ROIs as spheres centered on the statistical peaks reported in an  
17 activation likelihood estimation (ALE) meta-analysis of 16 fMRI studies (including 291  
18 subjects) involving theory-of-mind reasoning (van Veluw and Chance, 2014), in  
19 accordance with generally accepted guidelines in ROI analysis (Poldrack, 2007). The  
20 ROIs are shown in Figure 2. The peaks were located in six areas: the left TPJ (Montreal  
21 Neurological Institute [MNI]: -52, -56, 24), right TPJ (MNI: 55, -53, 24), left STS (MNI:  
22 -59, -26, -9), right STS (MNI: 59, -18, -17), MPFC (MNI: 1, 58, 19), and the precuneus  
23 (MNI: -3, -56, 37). The radius of the ROI spheres was 10 mm, corresponding to the

1 approximate volume (4,000 mm<sup>3</sup>) of the largest clusters (TPJ and MPFC) reported in van  
2 Veluw and Chance (2014). The same sphere radius was used for all ROIs.

3 The fMRI data from all participants were analyzed with the Statistical Parametric  
4 Mapping software (SPM12) (Wellcome Department of Cognitive Neurology, London,  
5 UK) (Friston et al., 1994). We first used a conventional general linear model (GLM) to  
6 estimate regression (beta) coefficients for each individual trial (i.e., 100 regressors),  
7 focusing on the 10-s story presentation phase of each trial. One regressor of no interest  
8 modelled the 4-s probe statement phase across conditions. Each regressor was modeled  
9 with a boxcar function and convolved with the standard SPM12 hemodynamic response  
10 function. In addition, ten run-specific regressors controlling for baseline differences  
11 between runs, and six motion regressors, were included. The trialwise beta coefficients  
12 for the endogenous and exogenous conditions (i.e., 80 beta maps) were then submitted to  
13 subsequent multivariate analyses (Haxby et al., 2001).

14 The MVPA was carried out using The Decoding Toolbox (TDT) version 3.999  
15 (Hebart et al., 2015) for SPM. For each subject and ROI, we used linear support vector  
16 machines (SVMs, with the fixed regularization parameter of  $C = 1$ ) to compute decoding  
17 accuracies. To ensure independent training and testing data sets, we used leave-one-run-  
18 out cross-validation approach. For each fold, the betas across all training runs were  
19 normalized relative the mean and standard deviation, and the same Z-transformation was  
20 applied to the betas in the left-out test run (Misaki et al., 2010). An SVM was then trained  
21 to discriminate activity patterns belonging to the endogenous or exogenous trials in nine  
22 runs, and then tested on the left-out run, repeated for all runs, resulting in a run-average  
23 decoding accuracy for each ROI and subject.

1 For statistical inference, the true group mean decoding accuracy was compared to  
2 a null distribution of group mean accuracies obtained from permutation testing. The same  
3 MVPA was repeated within each subject and ROI using permuted condition labels  
4 (10,000 iterations). A p value was computed as  $(1 + \text{the number of permuted group}$   
5  $\text{accuracy values} > \text{true value}) / (1 + \text{the total number of permutations})$ . To control for  
6 multiple comparisons across the six ROIs, we used the false discovery rate (FDR)  
7 correction (Benjamini and Hochberg, 1995). In addition, we also computed a bootstrap  
8 distribution around the true group mean accuracy by resampling individual-subject mean  
9 accuracies with replacement (10,000 iterations), from which a 95% confidence interval  
10 (CI) was derived (Nakagawa and Cuthill, 2007). A corrected p value  $< 0.05$  in  
11 combination with a 95% CI that does not cross chance level were interpreted as a  
12 significant decoding effect at the group level (Nakagawa and Cuthill, 2007).

13 In addition, as further exploratory statistics beyond the targeted hypotheses of this  
14 study, we used a whole-brain searchlight analysis (Kriegeskorte et al., 2006) to test for  
15 possible areas of decoding outside the ROIs. This searchlight analysis is described in the  
16 *Supplementary Information* (SI Text S1, Figures S1-S4, and Tables S1-S4).

17

### 18 *Testing prediction 2*

19 The purpose of the second analysis was to determine whether the brain encoded  
20 information concerning the type of agent (self versus other) present in the stories. The  
21 analysis methods were the same as for testing hypothesis 1, except that for regressors of  
22 interest we used the self-related and other-related trials, collapsed across the type of

1 attention (exogenous or endogenous). Just as for hypothesis 1, we tested the six defined  
2 ROIs within the theory-of-mind network.

3

#### 4 *Testing prediction 3*

5         The purpose of the third analysis was to test for an interaction between the two  
6 variables (endogenous versus exogenous, and self versus other). We used MVPA to test  
7 whether the decoding for the type of attention was significantly different between the  
8 self-related and the other-related stories. The analysis methods were similar to those used  
9 for testing hypothesis 1 and 2, except in the following ways. We computed two MVPA  
10 decoding results, the first for distinguishing endogenous-self from exogenous-self stories,  
11 the second for distinguishing endogenous-other from exogenous-other stories. We then  
12 computed the difference between the two decoding results ([endogenous-self versus  
13 exogenous-self] – [endogenous-other versus exogenous-other]) to create a decoding  
14 difference score. Just as for hypothesis 1 and 2, we tested the six defined ROIs within the  
15 theory-of-mind network(van Veluw and Chance, 2014).

16

#### 17 *Testing prediction 4*

18         The purpose of the fourth analysis was to confirm whether our story stimuli  
19 engaged social cognition and thereby recruited brain areas within the expected theory of  
20 mind network. The analysis was meant as an added control to check the validity of the  
21 paradigm. The analysis methods were similar to those used for testing hypothesis 1-3,  
22 except in the following ways. We computed four MVPA decoding results: endogenous-  
23 self versus nonsocial, endogenous-other versus nonsocial, exogenous-self versus

1 nonsocial, and exogenous-other versus nonsocial. (Because using MVPA to compare two  
2 conditions requires equal numbers of trials in both conditions, it was not possible to use a  
3 single analysis to compare all 80 social trials to the 20 nonsocial trials.) Each analysis  
4 represents a separate, alternative way to assess the social-versus-nonsocial decoding. Just  
5 as for hypothesis 1-3, we tested the six defined ROIs within the theory-of-mind network  
6 (van Veluw and Chance, 2014).

7

### 8 *Eye tracking analysis*

9 Eye movements were recorded via an MRI-compatible infrared eye tracker (SR  
10 Research EyeLink 1000 Plus), mounted just below the projector screen, sampling at 1000  
11 Hz. Before each scanning session, a calibration routine on five screen locations was used  
12 and repeated until the maximum error for any point was less than 1°. The obtained eye  
13 position data was cleaned of artifacts related to blink events and smoothed using a 20-ms  
14 moving average. We then built an SVM decoding model analogue to the cross-validation  
15 approach used for the fMRI data, but here based purely on eye tracking data, to test  
16 whether eye movement dynamics alone were sufficient to decode the conditions of  
17 interest (endogenous versus exogenous, and self versus other). In keeping with a previous  
18 study (Schneider et al., 2013), we organized the data in the following way. The part of the  
19 display within which the stimuli appeared was divided into an 8 x 4 grid of 32 equally  
20 sized squares. The grid covered the screen area within which the stories were presented  
21 (see red outline in Figure S5), and approximately corresponded to the locations of  
22 individual words (four lines, with eight words per line). For each trial, the proportion of  
23 time that the subject fixated within each square (32 features) and the saccades between



1 those regions ( $32 \times 32 = 1024$  features) was calculated. These 1056 features, representing  
2 information about both where people were looking as well as saccade dynamics, were  
3 then averaged across repetitions for each of the four main conditions within each of the  
4 10 runs, yielding one eye movement feature vector per condition per run (per subject).  
5 The feature vectors were submitted to an SVM classifier ( $C = 1$ ). Using a leave-one-run-  
6 out approach, the SVM model was trained on endogenous versus exogenous story types,  
7 and then tested in the left-out run. At the group level, the decoding accuracies were tested  
8 against chance level using t-tests. A similar analysis was then performed on the contrast  
9 between self-related stories versus other-related stories. The results showed that  
10 endogenous-versus-exogenous and self-versus-other story types could not be decoded  
11 significantly better than chance using the pattern of eye movement. See *Supplementary*  
12 *Information* (SI Text S2 and Figure S5) for the results of the eye-tracking analysis.

13

#### 14 *Data availability*

15 The data that support the findings of this study are available at  
16 <https://figshare.com/s/c3463d15bc78106a1b5c>.

17

## 18 **Results**

### 19 *Prediction 1*

20 We hypothesized that participants would encode the attentional state of the agents  
21 in the stories in enough detail to distinguish between endogenous and exogenous  
22 attention, even though the difference between the story types was extremely subtle – only  
23 a few words that very slightly altered the semantic meaning of the sentences. We made

1 the strong prediction that decoding would be found within the set of brain areas typically  
2 included in the theory-of-mind cortical network. In particular, based on prior studies  
3 (Igelström et al., 2016; Kelly et al., 2014), we anticipated that the decoding would be  
4 most evident in the TPJ. Figure 2 shows six ROIs within the theory-of-mind network,  
5 based on a meta-analysis of previous theory-of-mind studies (van Veluw and Chance,  
6 2014). Figure 3A shows the results (see Table 1 for numerical details). Decoding  
7 accuracy for endogenous versus exogenous stories was significantly above chance for the  
8 left TPJ, and the significance of the left TPJ decoding survived a multiple comparison  
9 correction for the six ROIs (mean decoding accuracy 52.9%, 95% CI 50.7 to 55.2,  
10  $p_{\text{uncorrected}}=0.0046$ ,  $p_{\text{FDR-corrected}}=0.0276$ ). The results confirm our central prediction. The  
11 left TPJ showed significant decoding of the attentional state – exogenous versus  
12 endogenous – of agents in a story. (See *Supplementary Information* for the results of an  
13 exploratory, brain-wide, searchlight analysis.)

14

## 15 *Prediction 2*

16 We hypothesized that participants would process the distinction between the two  
17 types of agent in the stories (self versus other). We made the strong prediction that  
18 decoding would be found within the same set of ROIs in the theory-of-mind cortical  
19 network. Figure 3B shows the results (see Table 1 for numerical details). Decoding  
20 accuracy for self versus other stories was significantly above chance, and survived a  
21 multiple comparisons correction, for the left TPJ (mean decoding accuracy 53.0%, 95%  
22 CI 50.1 to 55.6,  $p_{\text{uncorrected}}=0.0053$ ,  $p_{\text{FDR-corrected}}=0.0210$ ), left STS (mean decoding  
23 accuracy 52.3%, 95% CI 50.6 to 54.1,  $p_{\text{uncorrected}}=0.0204$ ,  $p_{\text{FDR-corrected}}=0.0306$ ), MPFC

1 (mean decoding accuracy 52.6%, 95% CI 50.5 to 55.0,  $p_{\text{uncorrected}}=0.0105$ ,  $p_{\text{FDR-}}$   
2  $\text{corrected}=0.0210$ ), and precuneus (mean decoding accuracy 52.7%, 95% CI 50.4 to 55.0,  
3  $p_{\text{uncorrected}}=0.0099$ ,  $p_{\text{FDR-corrected}}=0.0210$ ). These results confirm that the present paradigm,  
4 using stories that are subtly different from each other, can obtain social cognitive results  
5 that are consistent with previous findings.

6

### 7 *Prediction 3*

8 We hypothesized that areas in the theory-of-mind network would not only encode  
9 the distinction between endogenous and exogenous attention, but do so to a significantly  
10 different extent in self-related stories than in other-related stories. However, the results  
11 showed no significant interaction in any of the ROIs (Figure 3C and Table 1). Thus, we  
12 found no support for prediction 3.

13

### 14 *Prediction 4*

15 Finally, we asked whether the activity in the theory-of-mind network would  
16 distinguish between social stories and nonsocial stories. This final analysis served as a  
17 control to check the validity of the story stimuli and confirm that they engaged social  
18 cognition as expected. We expected a signal of much greater magnitude in this analysis  
19 than in the analyses described above. The reason is that, as noted above, the types of  
20 social stories differed from each other by only a few words, and were nearly identical in  
21 semantic content; thus any brain signal reflecting those differences is expected to be  
22 subtle. The distinction between social and nonsocial stories, however, was much greater

1 semantically, and therefore the evidence of decoding in the brain is expected to be of  
2 greater magnitude.

3 Figure 4 shows the results (see Table 2 for numerical details). The results are  
4 separated into six ROIs, and for each ROI, separated into four individual analyses,  
5 corresponding to each of the four main social conditions contrasted with the nonsocial  
6 control. Decoding accuracy was significantly greater than chance in almost all analyses  
7 across the six ROIs. The right STS showed the least consistent evidence of decoding. The  
8 TPJ bilaterally and the precuneus showed the most consistent evidence of decoding.  
9 These results show strong evidence of decoding of the social versus nonsocial stimuli in  
10 the known theory-of-mind, cortical network.

11

## 12 **Discussion**

13 This study analyzed brain activity while people read stories about agents attending  
14 to objects in the environment. We examined whether specific brain areas could decode  
15 information about the type of attention referenced in the story (exogenous versus  
16 endogenous), and about the type of agent in the story (whether the agent was the subject  
17 reading the story or a different person). We hypothesized that if the brain constructs a  
18 model of attentional state that is used in social cognition, then areas of the brain known to  
19 be involved in social cognition should be able to distinguish between the two types of  
20 attention, exogenous and endogenous, represented in the stories. Our main analysis  
21 confirmed the hypothesis: the left TPJ showed significant decoding of information about  
22 endogenous versus exogenous attention. The finding is, arguably, remarkable, given that  
23 the semantic and wording difference between the two story types is extremely subtle.

1           These results support a new and growing body of evidence that the human brain  
2 constructs a model of attention to aid in theory of mind (Guterstam et al., 2019;  
3 Guterstam and Graziano, 2020; Kelly et al., 2014; Pesquita et al., 2016; Vernet et al.,  
4 2019). The model includes information about attention that is deeper and more complex  
5 than just gaze direction or an identification of the attended object. At least one aspect of  
6 attention incorporated into the model appears to be the manner in which attention moves  
7 to an object: endogenously (internally directed) or exogenously (externally induced). The  
8 processing of the model appears to engage the theory-of-mind cortical network. The left  
9 TPJ showed the strongest decoding result. It is not clear why the left hemisphere showed  
10 stronger activity than the right in the present task. Social cognition tasks often activate  
11 the TPJ bilaterally, but typically engage the right TPJ more (Saxe and Wexler, 2005).  
12 One speculation is that some aspect of the present task, perhaps explicitly instructing  
13 people that the task was a test of reading comprehension, caused an emphasis on  
14 linguistic processing, biasing the activity toward the left hemisphere. Other explanations  
15 for the left-hemisphere bias may also be possible.

16           We were also able to analyze brain areas involved in self-versus-other encoding.  
17 We found evidence of self-versus-other encoding in the left TPJ, left STS, MPFC and  
18 precuneus. The MPFC and precuneus have been previously implicated in self processing  
19 (Northoff et al., 2006; Ochsner et al., 2004; Passingham et al., 2010; Qin and Northoff,  
20 2011; van Veluw and Chance, 2014), and the TPJ is consistently activated in fMRI  
21 studies involving self-recognition (van Veluw and Chance, 2014) and first-person  
22 perspective taking (Ionta et al., 2011). These results lend confidence to the present  
23 paradigm, showing that even the very subtle differences between our story stimuli were

1 able to reveal cortical results consistent with previous studies.

2         Contrary to our prediction 3, we found no evidence for an interaction between  
3 attention type and agent type decoding in the theory-of-mind ROIs. (As noted in the  
4 *Supplementary Information*, during an exploratory searchlight analysis, we also found no  
5 evidence of an interaction effect in any other brain area.) Although it is possible that our  
6 paradigm was simply not sensitive enough to detect subtle interaction effects, these  
7 results suggest that the brain encodes information about attention type in a similar  
8 manner in the self and in others.

9         Finally, significant decoding of the social-versus-nonsocial distinction was obtained  
10 across most of the theory-of-mind ROIs. This finding confirmed the validity of the  
11 paradigm, and was expected based on previous experiments of the theory-of-mind  
12 network (Gallagher et al., 2000; Saxe and Kanwisher, 2003; van Veluw and Chance,  
13 2014; Vogeley et al., 2001).

14         The use of a story-reading paradigm allowed us to systematically manipulate the  
15 kind of attention represented in the stimulus while keeping other experimental factors  
16 close to identical. The endogenous and exogenous story versions differed only with  
17 respect to a few key words specifying the type of attention, while the rest of the stories  
18 were semantically the same. To avoid cognitive bias or expectation effects, the probe task  
19 performed by the subjects concerned details about the spatial context or the objects in the  
20 stories, effectively distracting subjects from the description of attention. A post-scan  
21 questionnaire confirmed that none of the subjects came close to figuring out the purpose  
22 of the experiment (which they had been told was a “Reading Comprehension  
23 Experiment”). The finding of brain areas that significantly decoded the type of attention,

1 despite the distinction between endogenous and exogenous attention being subtle and  
2 task-irrelevant, indicates that the human brain automatically, and possibly also implicitly  
3 (Pesquita et al., 2016), constructs a model of an agents' attention that specifies at least  
4 some dynamic aspects of how that attention is moving around the scene.

5

## 6 **References**

- 7 Abraham A, Pedregosa F, Eickenberg M, Gervais P, Mueller A, Kossaifi J, Gramfort A,  
8 Thirion B, Varoquaux G. 2014. Machine learning for neuroimaging with scikit-  
9 learn. *Front Neuroinformatics* **8**:14.
- 10 Avants BB, Epstein CL, Grossman M, Gee JC. 2008. Symmetric diffeomorphic image  
11 registration with cross-correlation: evaluating automated labeling of elderly and  
12 neurodegenerative brain. *Med Image Anal* **12**:26–41.
- 13 Baron-Cohen S. 1997. *Mindblindness: An Essay on Autism and Theory of Mind*. MIT  
14 Press.
- 15 Benjamini Y, Hochberg Y. 1995. Controlling the false discovery rate: a practical and  
16 powerful approach to multiple testing. *J R Stat Soc Ser B Methodol* **57**:289–300.
- 17 Brainard DH. 1997. The psychophysics toolbox. *Spat Vis* **10**:433–436.
- 18 Calder AJ, Lawrence AD, Keane J, Scott SK, Owen AM, Christoffels I, Young AW.  
19 2002. Reading the mind from eye gaze. *Neuropsychologia* **40**:1129–1138.
- 20 Cox RW. 1996. AFNI: software for analysis and visualization of functional magnetic  
21 resonance neuroimages. *Comput Biomed Res* **29**:162–173.
- 22 Esteban O, Markiewicz CJ, Blair RW, Moodie CA, Isik AI, Erramuzpe A, Kent JD,  
23 Goncalves M, DuPre E, Snyder M, Oya H, Ghosh SS, Wright J, Durnez J,  
24 Poldrack RA, Gorgolewski KJ. 2019. fMRIPrep: a robust preprocessing pipeline  
25 for functional MRI. *Nat Methods* **16**:111–116.
- 26 Fletcher PC, Happé F, Frith U, Baker SC, Dolan RJ, Frackowiak RSJ, Frith CD. 1995.  
27 Other minds in the brain: a functional imaging study of “theory of mind” in story  
28 comprehension. *Cognition* **57**:109–128.
- 29 Fonov VS, Evans AC, McKinstry RC, Almli CR, Collins DL. 2009. Unbiased nonlinear  
30 average age-appropriate brain templates from birth to adulthood. *NeuroImage*  
31 **47**:S102.

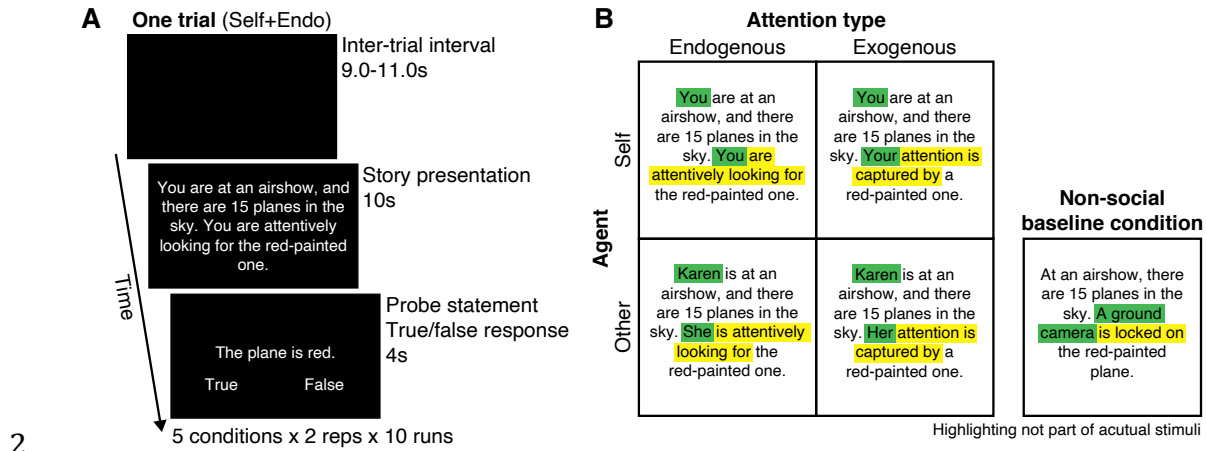
- 1 Friesen CK, Kingstone A. 1998. The eyes have it! Reflexive orienting is triggered by  
2 nonpredictive gaze. *Psychon Bull Rev* **5**:490–495.
- 3 Frischen A, Bayliss AP, Tipper SP. 2007. Gaze cueing of attention. *Psychol Bull*  
4 **133**:694–724.
- 5 Friston KJ, Holmes AP, Worsley KJ, Poline J-P, Frith CD, Frackowiak RS. 1994.  
6 Statistical parametric maps in functional imaging: a general linear approach. *Hum*  
7 *Brain Mapp* **2**:189–210.
- 8 Gallagher HL, Happé F, Brunswick N, Fletcher PC, Frith U, Frith CD. 2000. Reading the  
9 mind in cartoons and stories: an fMRI study of ‘theory of mind’ in verbal and  
10 nonverbal tasks. *Neuropsychologia* **38**:11–21.
- 11 Gorgolewski K, Burns CD, Madison C, Clark D, Halchenko YO, Waskom ML, Ghosh  
12 SS. 2011. Nipype: a flexible, lightweight and extensible neuroimaging data  
13 processing framework in python. *Front Neuroinformatics* **5**:13.
- 14 Graziano MS. 2019. Attributing awareness to others: The attention schema theory and its  
15 relationship to behavioural prediction. *J Conscious Stud* **26**:17–37.
- 16 Graziano MS. 2013. *Consciousness and the Social Brain*. Oxford University Press.
- 17 Graziano MS, Kastner S. 2011. Human consciousness and its relationship to social  
18 neuroscience: A novel hypothesis. *Cogn Neurosci* **2**:98–113.
- 19 Greve DN, Fischl B. 2009. Accurate and robust brain image alignment using boundary-  
20 based registration. *NeuroImage* **48**:63–72.
- 21 Guterstam A, Graziano MSA. 2020. Implied motion as a possible mechanism for  
22 encoding other people’s attention. *Prog Neurobiol* **190**:101797.
- 23 Guterstam A, Kean HH, Webb TW, Kean FS, Graziano MSA. 2019. Implicit model of  
24 other people’s visual attention as an invisible, force-carrying beam projecting  
25 from the eyes. *Proc Natl Acad Sci* **116**:328–333.
- 26 Guterstam A, Wilterson AI, Wachtell D, Graziano MSA. 2020. Other people’s gaze  
27 encoded as implied motion in the human brain. *Proc Natl Acad Sci* **117**:13162–  
28 13167.
- 29 Happé FG. 1994. An advanced test of theory of mind: Understanding of story characters’  
30 thoughts and feelings by able autistic, mentally handicapped, and normal children  
31 and adults. *J Autism Dev Disord* **24**:129–154.
- 32 Haxby JV, Gobbini MI, Furey ML, Ishai A, Schouten JL, Pietrini P. 2001. Distributed  
33 and overlapping representations of faces and objects in ventral temporal cortex.  
34 *Science* **293**:2425–2430.



- 1 Hebart MN, Gorgen K, Haynes J-D. 2015. The Decoding Toolbox (TDT): a versatile  
2 software package for multivariate analyses of functional imaging data. *Front*  
3 *Neuroinformatics* **8**.
- 4 Hoffman EA, Haxby JV. 2000. Distinct representations of eye gaze and identity in the  
5 distributed human neural system for face perception. *Nat Neurosci* **3**:80.
- 6 Igelstrom KM, Webb TW, Kelly YT, Graziano MSA. 2016. Topographical Organization  
7 of Attentional, Social, and Memory Processes in the Human Temporoparietal  
8 Cortex. *eNeuro* **3**.
- 9 Ionta S, Heydrich L, Lenggenhager B, Mouthon M, Fornari E, Chapis D, Gassert R,  
10 Blanke O. 2011. Multisensory Mechanisms in Temporo-Parietal Cortex Support  
11 Self-Location and First-Person Perspective. *Neuron* **70**:363–374.
- 12 Jenkinson M, Bannister P, Brady M, Smith S. 2002. Improved optimization for the robust  
13 and accurate linear registration and motion correction of brain images.  
14 *NeuroImage* **17**:825–841.
- 15 Kelly YT, Webb TW, Meier JD, Arcaro MJ, Graziano MS. 2014. Attributing awareness  
16 to oneself and to others. *Proc Natl Acad Sci* **111**:5012–5017.
- 17 Kobayashi H, Kohshima S. 1997. Unique morphology of the human eye. *Nature*  
18 **387**:767–768.
- 19 Kriegeskorte N, Goebel R, Bandettini P. 2006. Information-based functional brain  
20 mapping. *Proc Natl Acad Sci* **103**:3863–3868.
- 21 Logothetis NK, Pauls J, Augath M, Trinath T, Oeltermann A. 2001. Neurophysiological  
22 investigation of the basis of the fMRI signal. *Nature* **412**:150–157.
- 23 Misaki M, Kim Y, Bandettini PA, Kriegeskorte N. 2010. Comparison of multivariate  
24 classifiers and response normalizations for pattern-information fMRI.  
25 *NeuroImage* **53**:103–118.
- 26 Nakagawa S, Cuthill IC. 2007. Effect size, confidence interval and statistical  
27 significance: a practical guide for biologists. *Biol Rev* **82**:591–605.
- 28 Northoff G, Heinzel A, de Greck M, Bermpohl F, Dobrowolny H, Panksepp J. 2006.  
29 Self-referential processing in our brain—A meta-analysis of imaging studies on  
30 the self. *NeuroImage* **31**:440–457.
- 31 Ochsner KN, Knierim K, Ludlow DH, Hanelin J, Ramachandran T, Glover G, Mackey  
32 SC. 2004. Reflecting upon feelings: an fMRI study of neural systems supporting  
33 the attribution of emotion to self and other. *J Cogn Neurosci* **16**:1746–1772.
- 34 Passingham RE, Bengtsson SL, Lau HC. 2010. Medial frontal cortex: from self-generated  
35 action to reflection on one’s own performance. *Trends Cogn Sci* **14**:16–21.

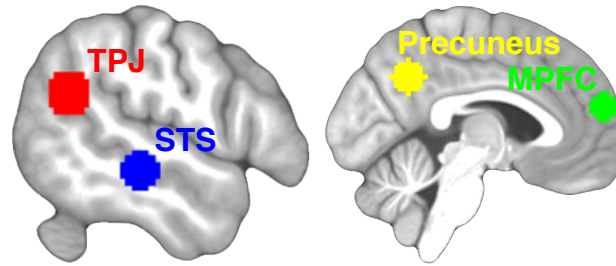
- 1 Perrett DI, Smith PAJ, Potter DD, Mistlin AJ, Head AS, Milner AD, Jeeves MA. 1985.  
2 Visual cells in the temporal cortex sensitive to face view and gaze direction. *Proc*  
3 *R Soc Lond B Biol Sci* **223**:293–317.
- 4 Pesquita A, Chapman CS, Enns JT. 2016. Humans are sensitive to attention control when  
5 predicting others' actions. *Proc Natl Acad Sci* **113**:8669–8674.
- 6 Poldrack RA. 2007. Region of interest analysis for fMRI. *Soc Cogn Affect Neurosci*  
7 **2**:67–70.
- 8 Puce A, Allison T, Bentin S, Gore JC, McCarthy G. 1998. Temporal Cortex Activation in  
9 Humans Viewing Eye and Mouth Movements. *J Neurosci* **18**:2188–2199.
- 10 Qin P, Northoff G. 2011. How is our self related to midline regions and the default-mode  
11 network? *NeuroImage*, Special Issue: Educational Neuroscience **57**:1221–1233.
- 12 Saxe R, Kanwisher N. 2003. People thinking about thinking people: the role of the  
13 temporo-parietal junction in “theory of mind.” *NeuroImage* **19**:1835–1842.
- 14 Saxe R, Wexler A. 2005. Making sense of another mind: The role of the right temporo-  
15 parietal junction. *Neuropsychologia* **43**:1391–1399.
- 16 Schneider B, Pao Y, Pea RD. 2013. Predicting Students' Learning Outcomes Using Eye-  
17 Tracking Data. *Learn Anal Knowl Conf*.
- 18 Tustison NJ, Avants BB, Cook PA, Zheng Y, Egan A, Yushkevich PA, Gee JC. 2010.  
19 N4ITK: improved N3 bias correction. *IEEE Trans Med Imaging* **29**:1310.
- 20 van Veluw SJ, Chance SA. 2014. Differentiating between self and others: an ALE meta-  
21 analysis of fMRI studies of self-recognition and theory of mind. *Brain Imaging*  
22 *Behav* **8**:24–38.
- 23 Vernet M, Japee S, Lokey S, Ahmed S, Zachariou V, Ungerleider LG. 2019. Endogenous  
24 visuospatial attention increases visual awareness independent of visual  
25 discrimination sensitivity. *Neuropsychologia*, Neural Routes to Awareness in  
26 Vision, Emotion and Action: A tribute to Larry Weiskrantz **128**:297–304.
- 27 Vogeley K, Bussfeld P, Newen A, Herrmann S, Happé F, Falkai P, Maier W, Shah NJ,  
28 Fink GR, Zilles K. 2001. Mind Reading: Neural Mechanisms of Theory of Mind  
29 and Self-Perspective. *NeuroImage* **14**:170–181.
- 30 Wicker B, Michel F, Henaff M-A, Decety J. 1998. Brain regions involved in the  
31 perception of gaze: a PET study. *NeuroImage* **8**:221–227.
- 32 Zhang Y, Brady M, Smith S. 2001. Segmentation of brain MR images through a hidden  
33 Markov random field model and the expectation-maximization algorithm. *IEEE*  
34 *Trans Med Imaging* **20**:45–57.

## 1 Figures and legends



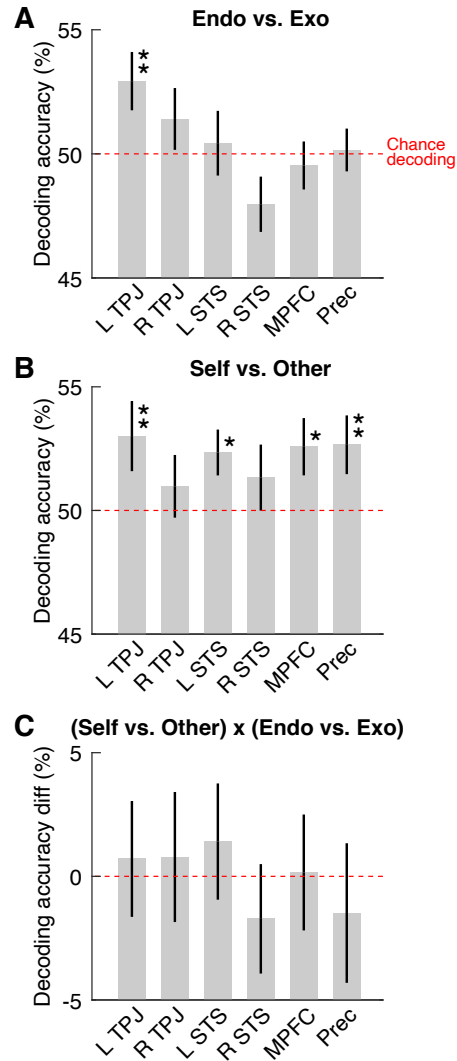
2

3 **Figure 1. Methods. A.** Schematic timeline of the fMRI design. In each trial, subjects  
4 were presented with a short story for 10 s, describing a scene in which an agent attended  
5 to an object in the environment. A probe statement was then shown for 4 s, relating to  
6 either the story's spatial context or object property, to which the subjects responded either  
7 true or false by button press. **B.** The agent in the story was either the subject him-/herself  
8 (self) or another person (other), and directed attention to the object endogenously  
9 (internally driven attention) or exogenously (stimulus-driven attention), yielding a  $2 \times 2$   
10 factorial design of attention type  $\times$  agent. We created 80 unique stories in four different  
11 versions, one for each condition. We made minimal changes to the wordings to keep the  
12 story versions as semantically similar as possible. Green highlighting indicates wording  
13 specifying agent, yellow highlighting indicates wording specifying attention type (colors  
14 not part of actual visual stimuli). For each story, each subject saw only one of the four  
15 versions (balanced across subjects). We also included a nonsocial baseline condition  
16 (twenty unique stories based on a subset of the 80 social stories) in which the agent was  
17 replaced by a non-human object.



1

2 **Figure 2. Regions of interest (ROIs).** Six ROIs were defined based on peaks reported in  
3 an activation likelihood estimation meta-analysis of 16 fMRI studies involving theory-of-  
4 mind reasoning (van Veluw and Chance, 2014). The ROIs consisted of 10-mm-radius  
5 spheres centered on peaks in the bilateral temporoparietal junction (TPJ) and superior  
6 temporal sulcus (STS), and two midline structures: the precuneus and medial prefrontal  
7 cortex (MPFC). Here, the TPJ and STS ROIs on the left side are shown.



1

2 **Figure 3. Decoding attention type, agent, and the interaction between them, in six**

3 **brain areas.** For definition of ROIs, see Figure 2. Each point shows mean decoding

4 accuracy. Error bars show SEM. Red horizontal line indicates chance level decoding.

5 Significance indicated by \* ( $p < 0.05$ ) and \*\* ( $p < 0.01$ ), based on permutation testing (all

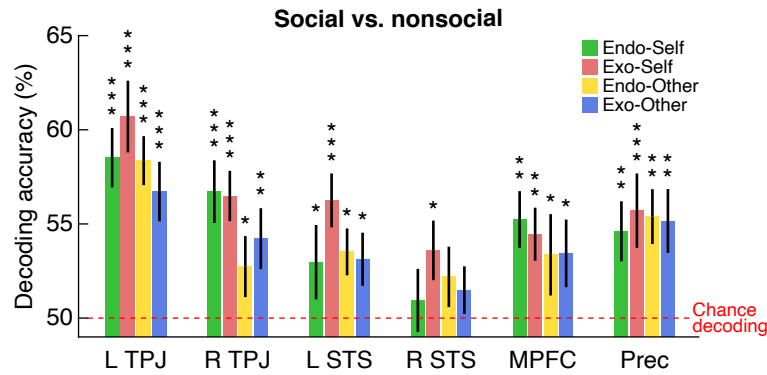
6 significant p values also survived correction for multiple comparisons across all six ROIs

7 [all corrected  $p$ s  $< 0.05$ ]). **A.** The ability of a classifier, trained on BOLD activity patterns

8 within each ROI, to decode endogenous (endo) versus exogenous (exo) attention. **B.**

9 Decoding accuracy for agent (self versus other). **C.** Decoding accuracy for the interaction

10 between type of attention and agent.



1

2 **Figure 4. Decoding social versus nonsocial stories.** The ability of a classifier, trained  
3 on BOLD activity patterns within each of the six ROIs, to decode each of the four social  
4 story conditions (endogenous-self, exogenous-self, endogenous-other, and exogenous-  
5 other) versus the nonsocial baseline. Each bar shows mean decoding accuracy, error bars  
6 show SEM, red horizontal line shows chance level decoding. Significance indicated by \*  
7 ( $p < 0.05$ ), \*\* ( $p < 0.01$ ), and \*\*\* ( $p < 0.001$ ) based on permutation testing (all but one of the  
8 significant p values also survived correction for multiple comparisons across all six ROIs;  
9 see Table 2 for numerical details).

## 1 Tables and legends

		L TPJ	R TPJ	L STS	R STS	MPFC	Precuneus
<b>Endo vs. Exo</b>	Mean accuracy	52.9%	51.4%	50.4%	48.0%	49.5%	50.2%
	95% CI	50.7 – 55.2	49.1 – 53.9	47.8 – 52.8	45.9 – 50.1	47.6 – 51.4	48.5 – 51.8
	P value	0.0046*	0.1148	0.3518	0.9547	0.6439	0.4428
<b>Self vs. Other</b>	Mean accuracy	53.0%	51.0%	52.3%	51.3%	52.6%	52.7%
	95% CI	50.1 – 55.6	48.5 – 53.4	50.6 – 54.1	48.9 – 54.1	50.5 – 55.0	50.4 – 55.0
	P value	0.0053*	0.1974	0.0204*	0.1241	0.0105*	0.0099*
<b>(Self vs. Other) × (Endo vs. Exo)</b>	Mean accuracy diff	1.6%	1.5%	2.0%	-3.0%	2.5%	0.6%
	95% CI	-2.7 – 6.3	-3.6 – 6.5	-2.7 – 6.3	-7.4 – 1.1	-2.3 – 6.7	-5.5 – 5.5
	P value	0.2430	0.2639	0.1967	0.8944	0.1414	0.3900

2 **Table 1. Decoding attention type, agent, and the interaction between the two, within**  
3 **the six ROIs.** For definition of ROIs, see Figure 2. Mean decoding accuracy (%), 95%  
4 confidence interval (based on bootstrap distribution), and p value (based on permutation  
5 testing) are shown for each of the six ROIs. Results shown for decoding endogenous  
6 (endo) versus exogenous (exo) attention type, self versus other agent type, and the  
7 interaction between the two variables. \*significant p values that survived correction for  
8 multiple comparisons across all six ROIs (corrected  $p < 0.05$ ).

		<b>L TPJ</b>	<b>R TPJ</b>	<b>L STS</b>	<b>R STS</b>	<b>MPFC</b>	<b>Precuneus</b>
<b>Endo-Self vs. nonsocial</b>	Mean accuracy	58.5%	56.7%	53.0%	50.9%	55.2%	54.6%
	95% CI	55.5 – 61.6	53.5 – 59.8	49.3 – 56.8	47.5 – 54.0	52.3 – 58.0	51.5 – 57.7
	P value	0.0001*	0.0001*	0.0338*	0.2703	0.0026*	0.0027*
<b>Exo-Self vs. nonsocial</b>	Mean accuracy	60.7%	56.5%	56.3%	53.6%	54.5%	55.7%
	95% CI	57.0 – 64.4	53.7 – 58.8	53.5 – 59.0	50.5 – 56.6	51.7 – 57.1	51.9 – 59.4
	P value	0.0001*	0.0002*	0.0001*	0.0179*	0.0044*	0.0001*
<b>Endo-Other vs. nonsocial</b>	Mean accuracy	58.4%	52.7%	53.5%	52.2%	53.4%	55.4%
	95% CI	55.9 – 60.9	49.6 – 55.9	51.1 – 55.9	48.8 – 55.0	49.5 – 57.7	52.6 – 58.2
	P value	0.0001*	0.0497	0.0147*	0.0969	0.0219*	0.0014*
<b>Exo-Other vs. nonsocial</b>	Mean accuracy	56.7%	54.2%	53.1%	51.5%	53.4%	55.2%
	95% CI	53.8 – 59.8	51.1 – 57.3	50.4 – 55.9	49.1 – 54.0	49.8 – 56.7	52.0 – 58.5
	P value	0.0002*	0.0065*	0.0319*	0.1901	0.0197*	0.0012*

1 **Table 2. Decoding social versus nonsocial stories within the six ROIs.** For definition  
2 of ROIs, see Figure 2. Mean decoding accuracy (%), 95% confidence interval (based on  
3 bootstrap distribution), and p value (based on permutation testing) are shown for each of  
4 the six ROIs. Results shown for each of four social story conditions (endogenous-self,  
5 exogenous-self, endogenous-other, and exogenous-other) versus the nonsocial baseline.  
6 \*significant p values that survived correction for multiple comparisons across all six  
7 ROIs (corrected  $p < 0.05$ ).

8



## 1 **Acknowledgments**

2 This work was supported by the Princeton Neuroscience Institute Innovation Fund. Arvid  
3 Guterstam was supported by the Wenner-Gren Foundation, the Sweden-America  
4 Foundation, and the Promobilia Foundation. The authors would like to thank Sam  
5 Nastase for valuable input regarding the multivoxel pattern analysis.

6

## 7 **Competing interests**

8 The authors declare no competing financial interests.

1 **Supplementary Information**

2

3 **The Supplementary Information (SI) includes:**

4 SI Text S1: Searchlight Analysis

5 SI Text S2: Eye-Tracking Decoding

6 Figures S1 to S5

7 Tables S1 to S4

8 SI Data S1: Story stimuli

## 1 **SI Text S1: Searchlight Analysis**

2 Our primary analysis, described in the main text, was a targeted testing of strong  
3 hypotheses within a set of defined RIOs. Such targeted testing is preferred because it puts  
4 strong hypotheses to a direct test, and because it is more sensitive, avoiding the statistical  
5 problems of a brain-wide multiple comparison correction. Given the extremely subtle  
6 differences between the story types in the present experiment, such a targeted and  
7 sensitive analysis was preferred. However, in addition to the targeted ROI analyses, we  
8 also performed an exploratory analysis to ask whether any meaningful decoding activity  
9 might be identified outside of the ROIs. We used a searchlight analysis (Kriegeskorte et  
10 al., 2006). The searchlight analysis is fundamentally different from the ROI analysis. It is  
11 not targeted to specific brain areas on the basis of predictions. Instead, it is a whole-brain  
12 analysis that is much more statistically conservative because of the brain-wide multiple  
13 comparisons. In general, one would not expect the searchlight analysis to align perfectly  
14 with the ROI analysis. Activations revealed in the more sensitive ROI analysis might not  
15 appear in the searchlight analysis. Instead, it is exploratory in nature, and its usefulness is  
16 that it may reveal clusters of strong decoding in unanticipated areas outside the ROIs.

17

### 18 *Endogenous-vs-exogenous searchlight analysis*

19 First, the brain was partitioned into overlapping voxel clusters of spherical shape  
20 (10-mm radius). In each of these clusters, a decoding accuracy was computed using the  
21 same model input, SVM parameters, and procedures as described for the ROI analysis.  
22 This process resulted in an endogenous-versus-exogenous decoding accuracy map for  
23 each subject, in which the value of each voxel represents the average proportion of

1 correctly classified trials relative to chance level (50%) based on the 10 mm sphere of  
2 tissue surrounding that voxel. The subject-wise decoding maps were then smoothed using  
3 a 3-mm full-width-half-maximum (FWHM) Gaussian kernel, and entered into a second-  
4 level analysis using SPM12. At the second-level, the whole-brain decoding maps were  
5 thresholded at  $p < 0.001$  (uncorrected for multiple comparisons). For statistical inference,  
6 we employed a cluster-level, whole-brain approach to find clusters that passed the  
7 threshold of  $p < 0.05$ , corrected for brain-wide multiple comparisons using the  
8 familywise error rate correction as implemented by SPM12. In a purely descriptive  
9 manner, we also report strong decoding activity, defined as clusters  $\geq 10$  voxels using the  
10 cluster-forming threshold of  $p < 0.001$  uncorrected, that did not survive correction at the  
11 whole-brain level (Table S1).

12 Clusters revealed in the searchlight analysis were projected onto orthogonal  
13 sections of the average structural scan generated from the 32 subjects for anatomical  
14 localization. The decoding clusters were also projected onto a 3D canonical brain surface  
15 using the software Surf Ice (University of South Carolina, McCausland Center for Brain  
16 Imaging). Figure S1 shows the brain-wide peak in decoding obtained with the searchlight  
17 method.

18 The searchlight analysis revealed no clusters that were brain-wide significant at the  
19 corrected  $p < 0.05$  threshold. When examining the voxel-wise  $p < 0.001$  threshold, four  
20 clusters ( $\geq 10$  voxels) were observed (see Table S1 for details). The strongest peak  
21 decoding ( $t = 4.21$ ) was located in the left posterior STS (within the TPJ; see Figure S1),  
22 thus consistent with the main, ROI analysis.

23

1 *Self-versus-other searchlight analysis*

2       The same method used for the exogenous-versus-endogenous searchlight analysis  
3 was used for the self-versus-other searchlight analysis. The analysis revealed four clusters  
4 that significantly decoded the self-versus-other distinction ( $p < 0.05$ ) after correcting for  
5 multiple comparisons using the whole brain as search space. The global decoding peak  
6 was located on the left angular gyrus (part of TPJ; see Figure S2) ( $t = 5.92$ ), which is  
7 compatible with our main ROI decoding results. All decoding clusters, including those  
8 passing the uncorrected threshold  $p < 0.001$ , are listed in Table S2.

9

10 *Endogenous-versus-exogenous X self-versus-other searchlight analysis*

11       The same method as in the previous sections was used for the interaction, or  
12 exogenous-versus-endogenous X self-versus-other, searchlight analysis. The analysis  
13 revealed no clusters that decoded the endogenous-versus-exogenous distinction  
14 significantly ( $p < 0.05$ , after correcting for multiple comparisons using the whole brain as  
15 search space) different in self-related compared to other-related stories. When examining  
16 the voxel-wise  $p < 0.001$  threshold, two clusters were observed (see Table S3 for details).  
17 The strongest peak decoding ( $t = 4.20$ ) was located in the left posterior STS (Figure S3).  
18 All decoding clusters, including those passing the uncorrected threshold  $p < 0.001$ , are  
19 listed in Table S3.

20

21 *Social-versus-nonsocial searchlight analysis*

22       For this analysis, we performed four separate whole-brain searchlight analyses,  
23 testing for decoding that distinguished social from nonsocial stories. Each of the four

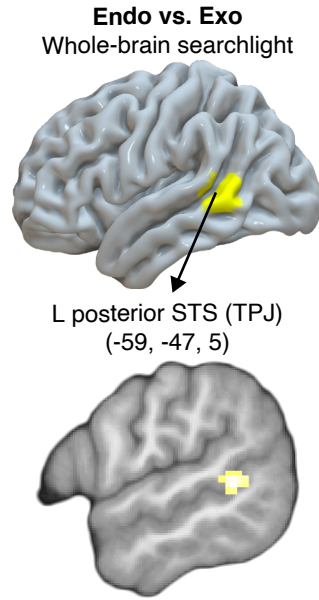
1 analyses was restricted to comparing a specific type of social story to the nonsocial  
2 control: endogenous-self versus nonsocial, exogenous-self versus nonsocial, endogenous-  
3 other versus nonsocial, and exogenous-other versus nonsocial. We first identified clusters  
4 that passed a threshold of  $p < 0.05$  corrected using the entire brain as search space. We  
5 then found brain areas of overlap, that were  $\geq 10$  voxels in size, between the decoding  
6 clusters obtained in the four different analyses. These areas of overlap represent brain  
7 regions that showed significant decoding for the social-versus-nonsocial comparison, in a  
8 consistent manner, across all types of social stimuli. Using this method, four brain areas  
9 were obtained (see Table S4 and Figure S4 for details).

10

11

## 1 **SI Text S2: Eye-Tracking Decoding**

2 To examine whether eye movement dynamics could explain our fMRI decoding  
3 results, we systematically organized eye position and saccade data into gaze pattern  
4 vectors, and submitted it to a decoding model analogue to the one used for the fMRI data  
5 (see Methods for details). Out of 32 subjects, 27 had usable eye tracking data available  
6 (five subjects had data with unacceptable levels of noise due to either the presence of  
7 glasses or a partially occluded pupil). The results showed that this classifier, based solely  
8 on eye tracking data, could not decode attention type (endogenous-versus-exogenous  
9 decoding accuracy 55.2%,  $p=0.102$ ) or agent (self-versus-other decoding accuracy  
10 52.4%,  $p=0.376$ ) significantly better than chance level. See Figure S5.

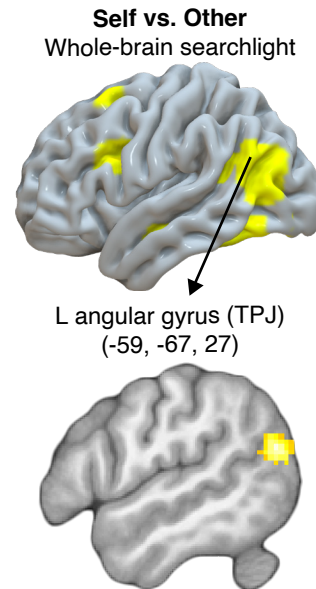


1

2 **Figure S1. Decoding attention type at the whole-brain level.** The cluster shown had  
3 the highest decoding accuracy in the whole-brain, searchlight analysis, for the  
4 endogenous-versus-exogenous comparison. See Table S1 for numerical details. Top:  
5 projected onto a 3D canonical brain surface. Bottom: projected onto a parasagittal section  
6 of the average structural scan generated from the 32 subjects for anatomical localization.  
7 For display purposes, the statistical threshold for the activation maps was set to  $p <$   
8 0.001, uncorrected.

9

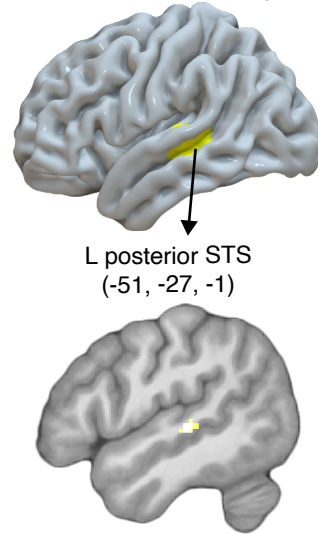




1

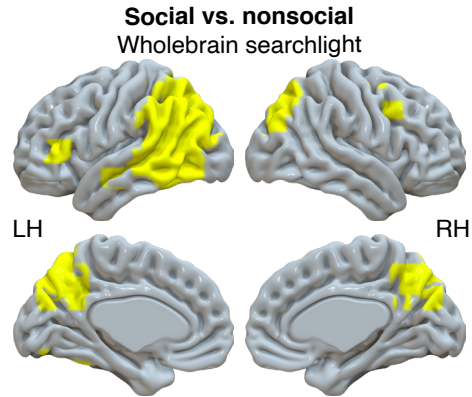
2 **Figure S2. Decoding agent type at the whole-brain level.** The cluster shown had the  
3 highest decoding accuracy in the whole-brain, searchlight analysis, for the self-versus-  
4 other comparison. See Table S2 for numerical details. Top: projected onto a 3D canonical  
5 brain surface. Bottom: projected onto a parasagittal section of the average structural scan  
6 generated from the 32 subjects for anatomical localization. For display purposes, the  
7 statistical threshold for the activation maps was set to  $p < 0.001$ , uncorrected.

**(Self vs. Other) x (Endo vs. Exo)**  
Whole-brain searchlight



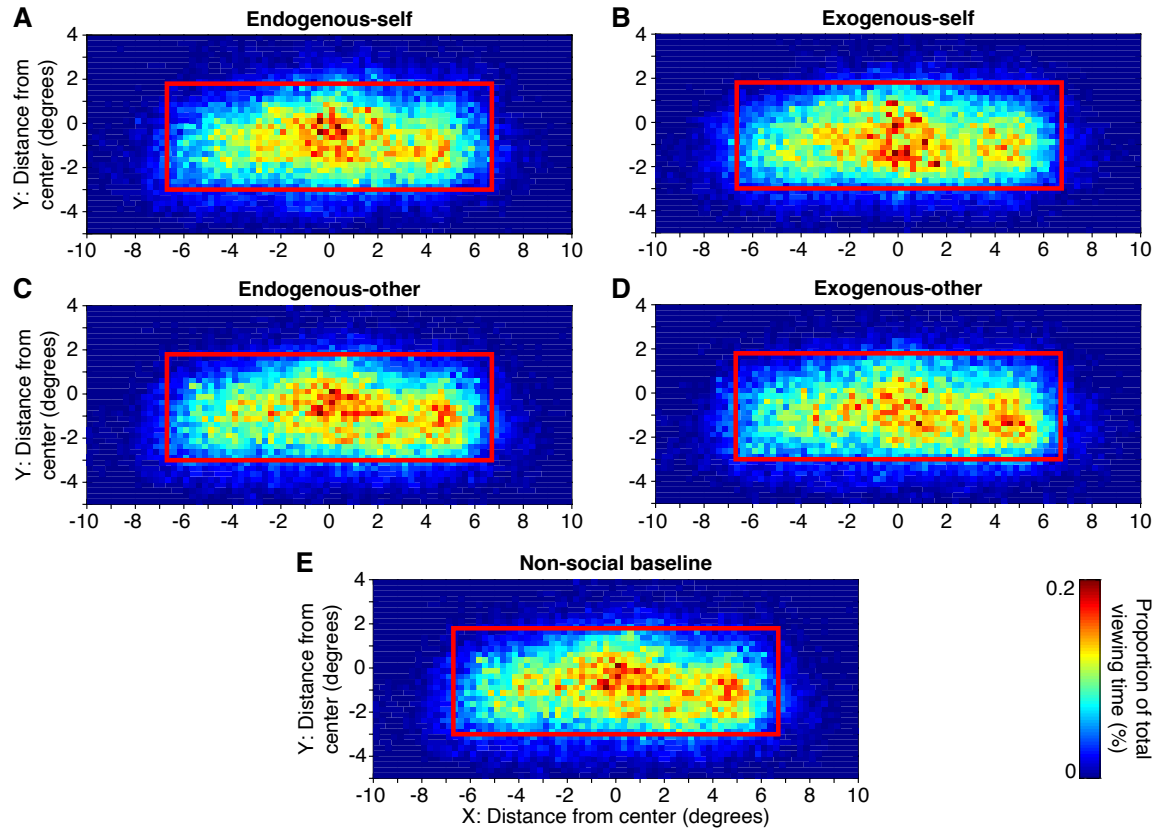
1

2 **Figure S3. Decoding attention-by-agent interaction at the whole-brain level.** The  
3 cluster shown had the highest decoding accuracy difference in the whole-brain,  
4 searchlight analysis, for the (endogenous-versus-exogenous)<sub>SELF</sub> versus (endogenous-  
5 versus-exogenous)<sub>OTHER</sub> comparison. See Table S3 for numerical details. Top: projected  
6 onto a 3D canonical brain surface. Bottom: projected onto a parasagittal section of the  
7 average structural scan generated from the 32 subjects for anatomical localization. For  
8 display purposes, the statistical threshold for the activation maps was set to  $p < 0.001$ ,  
9 uncorrected.



1

2 **Figure S4. Decoding social versus nonsocial stories at the whole-brain level.** The  
3 clusters shown are areas of overlap between four whole-brain, searchlight analyses, for  
4 the social-versus-nonsocial comparison (endogenous-self versus nonsocial, exogenous-  
5 self versus nonsocial, endogenous-other versus nonsocial, exogenous-other versus  
6 nonsocial). See Table S4 for numerical details. Projected onto a 3D canonical brain  
7 surface.



1

2 **Figure S5. Eye tracking results.** Subjects tended to fixate in a similar spatial pattern  
3 across the screen, and engage in similar saccade dynamics, regardless of the type of story  
4 presented. No significant ability to decode the story type based on the pattern of eye  
5 movement was obtained. The red rectangle in the heat maps in panels A-E indicates the  
6 area of the screen within which the story text appeared.

<b>Anatomical region</b>	<b>Peak MNI</b>	<b>Peak t</b>	<b>Cluster size</b>
L. posterior STS (TPJ)	-59, -47, 5	4.21	22
R. middle frontal gyrus	39, -7, 55	4.15	17
R. inferior precentral sulcus	42, 18, 15	4.10	37
L. lingual gyrus	-11, -80, -8	3.97	11

1 **Table S1. Decoding endogenous versus exogenous at the whole-brain level.** All  
2 clusters ( $\geq 10$  voxels) of decoding activity passing the voxelwise threshold of  $p < 0.001$   
3 (none of the clusters survived  $p < 0.05$  correction for multiple comparisons using the  
4 whole brain as search space).

Anatomical region	Peak MNI	Peak t	Cluster size	Corrected p value
L. angular gyrus (TPJ)	-59, -67, 27	5.92	211	<0.001
R. precuneus	4, -55, 30	5.91	57	-
L. inferior temporal gyrus	-41, -15, -18	5.49	131	0.004
L. inferior occipital gyrus	-31, -77, -8	5.26	54	-
R. middle cingulate cortex	2, -15, 40	4.24	47	-
R. lingual gyrus	24, -82, -8	4.82	121	0.006
L. fusiform gyrus	-39, -57, -23	4.81	114	0.008
R. calcarine sulcus	7, -85, 5	4.71	53	-
R. fusiform gyrus	27, -37, -21	4.32	34	-
R. inferior occipital gyrus	47, -80, 3	4.05	27	-
L. inferior frontal sulcus	-41, 13, 27	4.02	13	-
L. superior frontal gyrus	-29, 16, 65	3.98	12	-

1

2 **Table S2. Decoding agent type at the whole-brain level.** All clusters ( $\geq 10$  voxels)

3 decoding agent type (self versus other) activity at the threshold of  $p < 0.001$

4 (uncorrected). Corrected p values represent cluster-level correction using the whole brain

5 as search space.

<b>Anatomical region</b>	<b>Peak MNI</b>	<b>Peak t</b>	<b>Cluster size</b>
L. posterior STS	-51, -27, -1	4.20	13
R. inferior temporal gyrus	47, -22, -28	3.96	13

1 **Table S3. Decoding attention-by-agent interaction at the whole-brain level.** All  
2 clusters ( $\geq 10$  voxels) in which endogenous-versus-exogenous decoding was better in self-  
3 related compared to other-related stories at  $p < 0.001$  (uncorrected). (none of the clusters  
4 survived correction for multiple comparisons using the whole brain as search space).

<b>Anatomical region</b>	<b>Cluster size</b>
L. TPJ	1722
L. and R. precuneus	80
R. intraparietal sulcus	137
L. inferior frontal sulcus	22

1 **Table S4. Decoding social versus nonsocial stories at the whole-brain level.** Clusters  
2 ( $\geq 10$  voxels) decoding social versus nonsocial stories significantly better than chance.  
3 The listed clusters represent the overlap of significant clusters ( $p < 0.05$ , corrected using a  
4 cluster-defining uncorrected threshold of  $p < 0.001$  and the entire brain as search space)  
5 across four separate whole-brain searchlight analyses: endogenous-self versus nonsocial,  
6 exogenous-self versus nonsocial, endogenous-other versus nonsocial, and exogenous-  
7 other versus nonsocial.



- 1 **SI Data S1. Story stimuli.** List of all of the stories and story versions (endogenous-self,
- 2 exogenous-self, endogenous-other, and exogenous-other) presented to the participants
- 3 during the experiment.