



ELSEVIER

NeuroImage

www.elsevier.com/locate/ynimg
NeuroImage xx (2008) xxx–xxx

The neural basis of belief encoding and integration in moral judgment

Liane Young^{a,b,*} and Rebecca Saxe^b

^aDepartment of Psychology, Harvard University, 33 Kirkland Street, Cambridge, MA 02138, USA

^bDepartment of Brain and Cognitive Sciences, Massachusetts Institute of Technology, 43 Vassar Street, Cambridge, MA 02139, USA

Received 6 November 2007; revised 13 January 2008; accepted 19 January 2008

Moral judgment in the mature state depends on “theory of mind”, or the capacity to attribute mental states (e.g., beliefs, desires, and intentions) to moral agents. The current study uses functional magnetic resonance imaging (fMRI) to investigate the cognitive processes for belief attribution in moral judgment. Participants read vignettes in a 2×2×2 design: protagonists produced either a negative or neutral outcome, based on the belief that they were causing the negative outcome or the neutral outcome; presentation of belief information either preceded or followed outcome information. In each case, participants judged the moral permissibility of the action. The results indicate that while the medial prefrontal cortex is recruited for processing belief valence, the temporo-parietal junction and precuneus are recruited for processing beliefs in moral judgment via two distinct component processes: (1) encoding beliefs and (2) integrating beliefs with other relevant features of the action (e.g., the outcome) for moral judgment.

© 2008 Published by Elsevier Inc.

Keywords: Morality; Theory of mind; Belief attribution; fMRI; Temporo-parietal junction; Precuneus; Medial prefrontal cortex

Introduction

One key cognitive input to moral judgment is “theory of mind” or the capacity to attribute mental states, such as beliefs, desires, and intentions, to moral agents (e.g., Baird and Astington, 2004; Borg et al., 2006; Cushman et al., 2006; Knobe, 2005; Mikhail, 2007; Young et al., 2007). Adults judge intentional harms to be morally worse than the same harms brought about accidentally or unknowingly. In the current study, we investigate the neural evidence for multiple distinct cognitive processes underlying theory of mind in moral judgment.

The neural basis of theory of mind has been investigated in recent functional magnetic resonance imaging (fMRI) studies. These studies reveal a consistent group of brain regions for “theory of mind” in nonmoral contexts: the medial prefrontal cortex, right and left temporo-parietal junction, and precuneus (Ciaramidaro et al., 2007; Fletcher et al., 1995; Gallagher et al., 2000; Gobbini et al., 2007; Ruby and Decety, 2003; Saxe and Kanwisher, 2003; Vogeley et al., 2001). Of these regions, the right temporo-parietal junction (RTPJ) in particular appears to be selective for belief attribution (Aichorn et al., in press; Fletcher et al., 1995; Gallagher et al., 2000; Gobbini et al., 2007; Perner et al., 2006; Saxe and Wexler, 2005; Sommer et al., 2007). For example, its response is high when subjects read stories that describe a character’s thoughts and beliefs but low during stories containing other socially relevant information (e.g., a character’s physical appearance, cultural background, or even internal subjective sensations such as hunger or fatigue; Saxe and Powell, 2006).

A recent fMRI study showed that these same brain regions are recruited for moral judgment, particularly, judgment of intentional and unintentional harms and non-harms (Young et al., 2007). These brain regions showed significant activation above baseline for all conditions of moral judgment but were modulated by an interaction between mental state and outcome factors. In the current study, we sought to refine our characterization of the role of these brain regions. Evidence from developmental psychology suggests that the acquisition of the theory of mind skills required for mature moral judgment is marked by multiple distinct cognitive achievements. We investigated whether these different developmental stages correspond to distinct functional profiles in the adult brain.

The classic task for assessing a child’s ability to reason about the mental states of others is the false belief task (for a review, see Flavell, 1999; Wellman et al., 2001). In its standard version, known as the “object transfer” problem, the child is told a story in which a character’s belief about the location of a target object becomes false when the object is moved without the character’s knowledge. Generating the correct answer requires the child to pay attention to the character’s belief, not just to the true location of the object. While the precise age of success varies between children and between versions of the task, in general, children younger than 3 or 79

* Corresponding author. Department of Psychology, Harvard University, 33 Kirkland Street, Cambridge, MA 02138, USA. Fax: +1 617 258 8654.

E-mail address: lyoung@fas.harvard.edu (L. Young).

Available online on ScienceDirect (www.sciencedirect.com).

80 4 years old cannot verbalize correct answers to false belief
81 problems (but see Onishi and Baillargeon, 2005). By the time they
82 are five, children reliably pass the false belief test.

83 This capacity appears to precede rather than to coincide with the
84 capacity to use belief information in the context of moral judgment.
85 Five year olds can make moral distinctions based on mental state
86 distinctions only when consequences are held constant (Karniol,
87 1978; Nelson Le Gall, 1985; Nunez and Harris, 1998; Siegel and
88 Peterson, 1998; Wellman et al., 1979). Even though they can
89 represent beliefs, these children continue to base their moral
90 judgments primarily on the action's consequences rather than the
91 actor's beliefs, when these two factors conflict (Hebble, 1971; Piaget,
92 1965/1932; Shultz et al., 1986; Yuill, 1984; Yuill and Perner, 1988;
93 Zelazo et al., 1996). For example, five year olds judge that an agent
94 who intends to direct a traveler to the right location but accidentally
95 misdirects him is worse than another agent who intends to misdirect a
96 traveler but accidentally directs him to the right place (Piaget, 1965/
97 1932). Only later are children able to generate adult-like judgments
98 of these scenarios, which continue to take consequences into account
99 (Cushman, personal communication) but additionally depend
100 substantially on beliefs (Baird and Astington, 2004; Baird and
101 Moses, 2001; Darley and Zanna, 1982; Fincham and Jaspers, 1979;
102 Karniol, 1978; Shultz et al., 1986; Yuill, 1984) thereby requiring true
103 integration of information about consequences and beliefs (Grue-
104 neich, 1982; Weiner, 1995; Zelazo et al., 1996).

105 Based on this evidence from developmental psychology, we
106 propose a distinction between two separate component processes of
107 belief attribution in moral judgment: encoding and integration.
108 Encoding consists of forming an initial representation of the
109 protagonist's belief. Integration, by contrast, consists of using the
110 belief for moral judgment in flexible combination with relevant
111 outcome information. On this analysis, five-year-old children are
112 capable of encoding beliefs (e.g., in the false belief task), but they
113 cannot fully integrate beliefs with outcomes in the service of moral
114 judgment. Here we investigate the neural evidence for these
115 cognitive processes in the adult brain. We suggest that the brain
116 regions for encoding should be (1) recruited when belief information
117 is first presented and (2) recruited selectively for belief information
118 over non-belief information. As such, the response at encoding
119 should be stimulus-bound, that is, modulated by whether the current
120 stimulus being processed contains belief content. Brain regions for
121 integration should be (1) recruited once morally relevant non-belief
122 information (e.g., outcome) is available and (2) show a functional
123 profile reflecting the interaction between belief and outcome. The
124 response at integration should therefore reflect the use of prior belief
125 information in constructing a moral judgment and the influence of
126 outcome information on belief processing.

127 In the current study, participants read vignettes in a $2 \times 2 \times 2$
128 design (Fig. 1): protagonists produced either a negative outcome or a
129 neutral outcome, based on the belief that they were causing the
130 negative outcome ("negative" belief) or the neutral outcome
131 ("neutral" belief); belief information could be presented either
132 before or after information foreshadowing the outcome. A
133 protagonist with a negative belief who produced a negative outcome
134 did so knowingly, while a protagonist with a negative belief who
135 produced a neutral outcome did so unknowingly or accidentally,
136 based on a false belief. In each case, participants judged the moral
137 permissibility of the protagonist's action. This design allowed us to
138 address the following questions with respect to theory of mind in
139 mature moral judgment: (1) Is there neural evidence for encoding
140 and integration as distinct processes? (2) Are brain regions pre-

viously implicated in belief attribution in nonmoral contexts spe- 141
cifically involved in belief encoding and/or belief integration? (3) If 142
so, are encoding and integration accomplished by the same or dif- 143
ferent subsets of these regions? 144

145 Materials and methods

Seventeen naive right-handed subjects (Harvard College under- 146
graduates, aged 18–22 years, six women) participated in the func- 147
tional MRI study for payment. All subjects were native English 148
speakers, had normal or corrected-to-normal vision and gave 149
written informed consent in accordance with the requirements of 150
the internal review board at MIT. Subjects were scanned at 3-T (at 151
the MIT scanning facility in Cambridge, Massachusetts) using 152
twenty-six 4-mm-thick near-axial slices covering the whole brain. 153
Standard echoplanar imaging procedures were used (TR=2 s, 154
TE=40 ms, flip angle 90°). 155

Stimuli consisted of eight variations of 48 scenarios for a total of 156
384 stories with an average of 86 words per story (see Sup- 157
plementary data for full text of scenarios). A $2 \times 2 \times 2$ design was 158
used for each scenario: (i) protagonists produced either a negative 159
outcome (harm to a person) or a neutral outcome (no harm); (ii) 160
protagonists held the belief that they were causing a negative 161
outcome ("negative" belief) or a neutral outcome ("neutral" belief); 162
(iii) either belief information or information foreshadowing the 163
outcome was presented first. Stories were presented in four cu- 164
mulative segments, each presented for 6 s, for a total presentation 165
time of 24 s per story: 166

- (1) Background: information to set the scene (identical across 167
all conditions) 168
- (2 or 3) Foreshadow: information foreshadowing the outcome (nega- 169
tive or neutral) 170
- (2 or 3) Belief: the protagonist's belief about the situation (negative 171
or neutral) 172
- (4) Outcome: the protagonist's action and its outcome (nega- 173
tive or neutral) 174
175

For example, as in the scenario in Fig. 1, the identification of the 176
white powder by the coffee as poison rather than sugar foreshadows 177
a person's death by poison. In every story used in this experiment, 178
when something is wrong at this stage (e.g., poison in place of sugar, 179
drowning swimmer, asthma attack), the protagonist's action or 180
inaction results in a negative outcome (someone's death). Each 181
possible belief was true for one outcome and false for the other 182
outcome. Stories were presented and then removed from the screen 183
and replaced with a question about the moral nature of the action. 184
Participants were asked to make judgments on a scale of 1 185
(forbidden) to 3 (permissible), using a button press. Three buttons 186
were used due to the malfunction of a fourth button on the scanner- 187
safe response apparatus. The question remained on the screen for 4 s. 188

Subjects saw one variation of each scenario, for a total of 48 189
stories. Stories were presented in a pseudorandom order, the order of 190
conditions counterbalanced across runs and across subjects, thereby 191
ensuring that no condition was immediately repeated. Eight stories 192
were presented in each 5.6 min run; the total experiment, involving 193
six runs, lasted 33.6 min. Fixation blocks of 14 s were interleaved 194
between each story. The text of the stories was presented in a white 195
24-point font on a black background. Stories were projected onto a 196
screen via Matlab 5.0 running on an Apple G4 laptop. 197

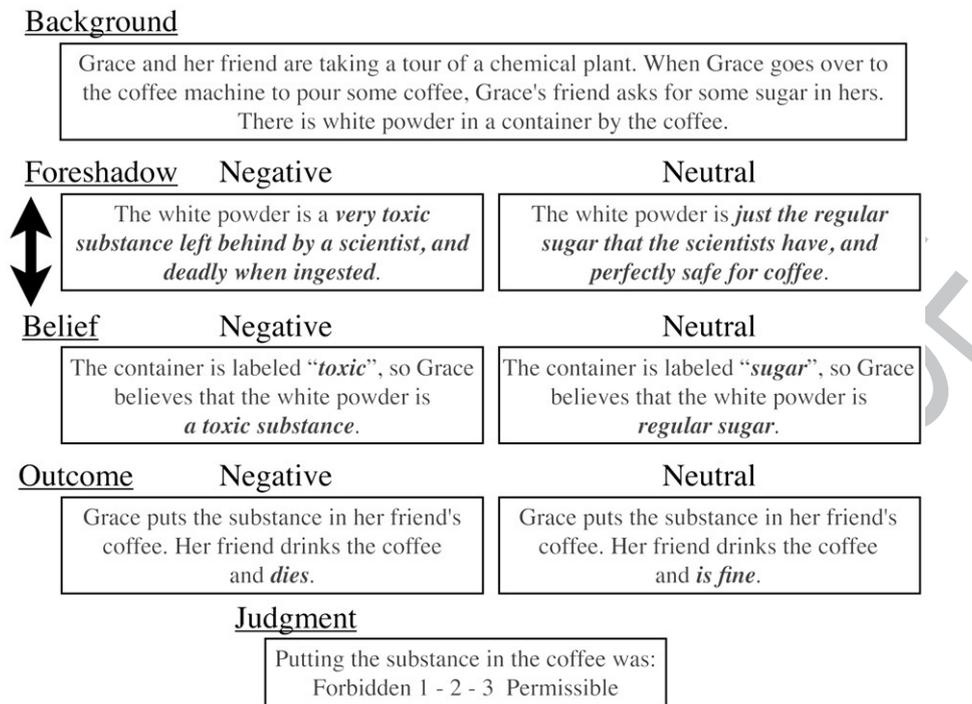


Fig. 1. Schematic representation of sample scenario. Changes between conditions are highlighted in bold italics. "Foreshadow" information foreshadows whether the action will result in a negative or neutral outcome. "Belief" information states whether the protagonist holds a belief that he or she is in a negative situation and that action will result in a negative outcome (negative belief) or a belief that he or she is in a neutral situation and that action will result in a neutral outcome (neutral belief). During belief-first trials, belief information was presented first and foreshadow information was presented second. During foreshadow-first trials, foreshadow information was presented first and belief information was presented second. Sentences corresponding to each category were presented in 6 s blocks.

198 In the same scan session, subjects participated in four runs of a
199 localizer experiment, contrasting stories that required inferences
200 about a character's beliefs (belief condition) with stories that
201 required inferences about a physical representation, i.e. a photo that
202 has become outdated (photo condition). Stimuli and story presenta-
203 tion were exactly as described in Saxe and Kanwisher (2003),
204 Experiment 2.

205 *fMRI data analysis*

206 MRI data were analyzed using SPM2 (<http://www.fil.ion.ucl.ac.uk/spm>) and custom software. Each subject's data were motion
207 corrected and then normalized onto a common brain space (the
208 Montreal Neurological Institute, MNI, template). Data were then
209 smoothed using a Gaussian filter (full width half maximum = 5 mm),
210 and high-pass filtered during analysis. A slow event-related design
211 was used and modeled using a boxcar regressor. An event was
212 defined as a single story (30 s); the event onset was defined by the
213 onset of text on the screen. The timing of the four story components
214 was constant for every story; thus, independent parameter estimates
215 were not created for each component. Components were separated
216 by the time of response, accounting for the hemodynamic lag.

217 Both whole-brain and tailored regions of interest (ROI) analyses
218 were conducted. Six ROIs were defined for each subject individually
219 based on a whole-brain analysis of a localizer contrast, and defined
220 as contiguous voxels that were significantly more active ($p < 0.001$,
221 uncorrected) while the subject read belief stories, as compared
222 with photo stories: RTPJ, LTPJ, PC, dMPFC, mMPFC, and vMPFC.
223 All peak voxels are reported in Montreal Neurological Institute
224 coordinates.
225

226 The responses of these regions of interest were then mea-
227 sured while subjects read stories from the current experiment.
228 Within the ROI, the average percent signal change (PSC) relative
229 to rest baseline ($PSC = 100 \times \text{raw BOLD magnitude for (condition}$
230 $- \text{fixation}) / \text{raw BOLD magnitude for fixation}$) was calculated
231 for each condition at each time point (averaging across all vo-
232 xels in the ROI and all blocks of the same condition). PSC during
233 story presentation (adjusted for hemodynamic lag) in each of
234 the ROIs was compared across experimental conditions. Because
235 the data defining the ROIs were independent from the data used
236 in the repeated measures statistics, Type I errors were drastically
237 reduced.

238 **Results and discussion**

239 *Behavioral results*

240 Subjects evaluated the moral status of protagonists' actions using
241 three buttons associated with a scale from completely forbidden
242 (1) to completely permissible (3). To determine the effects of belief
243 and outcome and order, we used a $2 \times 2 \times 2$ (outcome [negative vs.
244 neutral] by belief ["negative" vs. "neutral"] by order [belief-first vs.
245 foreshadow-first]) repeated measures ANOVA. Actions performed
246 by protagonists with "negative" beliefs were judged to be less per-
247 missible than actions performed by protagonists with "neutral" be-
248 liefs (negative: 1.2, neutral: 2.2; $F(1,11) = 69.7$, $p = 4.4 \times 10^{-6}$, partial
249 $\eta^2 = 0.86$). Actions resulting in negative outcomes were judged to
250 be less permissible than actions resulting in neutral outcomes (nega-
251 tive: 2.1, neutral: 2.5; $F(1,11) = 20.4$, $p = 0.001$, partial $\eta^2 = 0.65$). No
252 other main effect or interaction achieved significance. The same

253 $2 \times 2 \times 2$ repeated measures ANOVA was performed for reaction
254 time, yielding no significant main effects or interactions.

255 *fMRI results: localizer task*

256 To define regions implicated in belief attribution, stories that
257 required inferences about a character's beliefs (belief condition)
258 were contrasted with stories that required inferences about a
259 physical representation such as an outdated photograph (photo
260 condition). A whole-brain random effects analysis of the data
261 replicated results of previous studies using the same task (Saxe and
262 Kanwisher, 2003; Saxe and Wexler, 2005), revealing a higher
263 BOLD response during belief, as compared to photo stories, in the
264 RTPJ, LTPJ, dorsal (d), middle (m), and ventral (v) MPFC,
265 precuneus (PC), right temporal pole, and right anterior superior
266 temporal sulcus ($p < 0.001$, uncorrected, $k > 10$). Regions of interest
267 (ROIs) were identified in individual subjects (Table 1) at the same
268 threshold: RTPJ (15/17 subjects), PC (17/17), LTPJ (16/17),
269 dMPFC (14/17), mMPFC (12/17), and vMPFC (10/17).

270 *fMRI results: moral judgment task*

271 The average percent signal change (PSC) from rest in each
272 region of interest was calculated for each of three time intervals:

273 Time 1 (10–14 s): belief (belief-first trials) or foreshadow (fore-
274 shadow-first trials)

275 Time 2 (16–20 s): foreshadow (belief-first trials) or belief (fore-
276 shadow-first trials)

277 Time 3 (22–26 s): information about the protagonist's action

278
279 Times 1 and 2 represent the time during which the *encoding* of
280 the belief may occur; belief information is being presented for the
281 first time, and information relevant for moral judgment is incom-
282 plete. Time 3 represents the time during which the *integration* of
283 the belief may occur; no new belief information is added, but prior
284 belief information may be integrated with information about the
285 protagonist's action and the actual outcome in the construction of a
286 moral judgment.

287 *Encoding*

288 The PSCs for the earlier times, times 1 and 2 (encoding), during
289 which belief and foreshadow information were initially presented,
290 were analyzed using a $2 \times 2 \times 2 \times 2$ (time [1 vs. 2] by outcome

t1.1 Table 1

t1.2 Localizer experiment results

t1.3 ROI	Individual ROIs			Whole-brain contrast		
	x	y	z	x	y	z
t1.5 RTPJ	56	-56	22	56	-54	28
t1.6 PC	-1	-58	39	-2	-60	40
t1.7 LTPJ	-50	-63	26	-52	-58	26
t1.8 dMPFC	-2	58	29	2	60	28
t1.9 mMPFC	1	59	15	-4	56	8
t1.10 vMPFC	1	55	-7	0	54	-8

Average peak voxels for ROIs in Montreal Neurological Institute coordinates.
The "Individual ROIs" columns show the average peak voxels for individual
subjects' ROIs. The "Whole-brain contrast" columns show the peak voxel in
the same regions in the whole-brain random effects group analysis.

t1.11

[negative vs. neutral] by belief ["negative" vs. "neutral"] by order
[belief-first vs. foreshadow-first]) repeated measures ANOVA. 292

(1) RTPJ: A significant time by order interaction was observed in
the RTPJ ($F(1,14)=8.0$, $p=0.01$): the average response was
higher at time 1 when belief information (mean PSC: 0.41) was
presented at time 1, than when foreshadow information (mean
PSC: 0.35) was presented at time 1; and higher at time 2 when
belief (mean PSC: 0.54) was presented at time 2 than when
foreshadow (mean PSC: 0.44) was presented at time 2. Planned
comparisons at times 1 and 2 did not yield significant
differences between belief and foreshadow, though averaging
over times 1 and 2 revealed a greater response for belief than
foreshadow (mean belief PSC: 0.47; mean foreshadow PSC: 0.40;
 $t(14)=2.82$, $p=0.01$). The PSC in the RTPJ therefore
appeared to track with whether the stimulus being presented
contained belief information or not (Fig. 2, top panel; Table 2).
However, the response during the encoding of the belief did not
depend on the content or "valence" of the belief. At the time
that the belief was presented, there was no difference between
the responses to "negative" versus "neutral" belief (belief at
time 1: "negative": 0.39, "neutral": 0.41, $t(14)=0.32$, $p=0.76$;
belief at time 2: "negative": 0.55, "neutral": 0.53, $t(14)=-0.32$,
 $p=0.75$). There were also no main effects of negative versus
neutral foreshadow during encoding.

(2) PC and LTPJ: The PC and LTPJ showed a similar though less
selective profile at encoding (Tables 2 and 3, Supplementary
Fig. 1). A time (1 vs. 2) by order (belief-first vs. foreshadow-
first) interaction was observed in both the PC ($F(1,16)=7.4$,
 $p=0.02$) and the LTPJ ($F(1,15)=5.0$, $p=0.04$). That is, the
response in the PC was higher at time 1 when belief infor-
mation was presented at time 1 and higher at time 2 when
belief information was presented at time 2, suggesting that the
response in the PC during encoding, like the RTPJ, is driven by
the stimulus—whether the stimulus contains belief informa-
tion. Like the RTPJ, planned comparisons for the PC at times 1
and 2 did not yield significant differences between belief and
foreshadow, though averaging over times 1 and 2 revealed a
greater response for belief than foreshadow (mean belief PSC:
0.07; mean foreshadow PSC: 0.001; $t(16)=2.71$, $p=0.02$). The
interaction was less selective in the LTPJ: belief versus
foreshadow was discriminated at time 2 but not at time 1.

(3) MPFC: Regions in the MPFC showed a different pattern from
the RTPJ, PC, and LTPJ. There was no evidence that any region
of the MPFC was recruited for belief encoding (Fig. 2, bottom
panel). No significant main effects or interactions were found
during times 1 and 2 in the dMPFC, mMPFC, or vMPFC. To
determine whether the profile found during encoding (times 1
and 2) for the RTPJ (e.g., time by order interaction) was
significantly different from the profile found for regions in the
MPFC, a $2 \times 2 \times 2$ repeated measures ANOVA was conducted
for every pair of regions that included the RTPJ and one region
in the MPFC. The predicted time by order by region interactions
were significant ($p < 0.05$) in all pairs.

Integration

The PSC for time 3 (integration) was analyzed using a $2 \times 2 \times 2$
(outcome [negative vs. neutral] by belief ["negative" vs. "neutral"]
by order [belief-first vs. foreshadow-first]) repeated measures

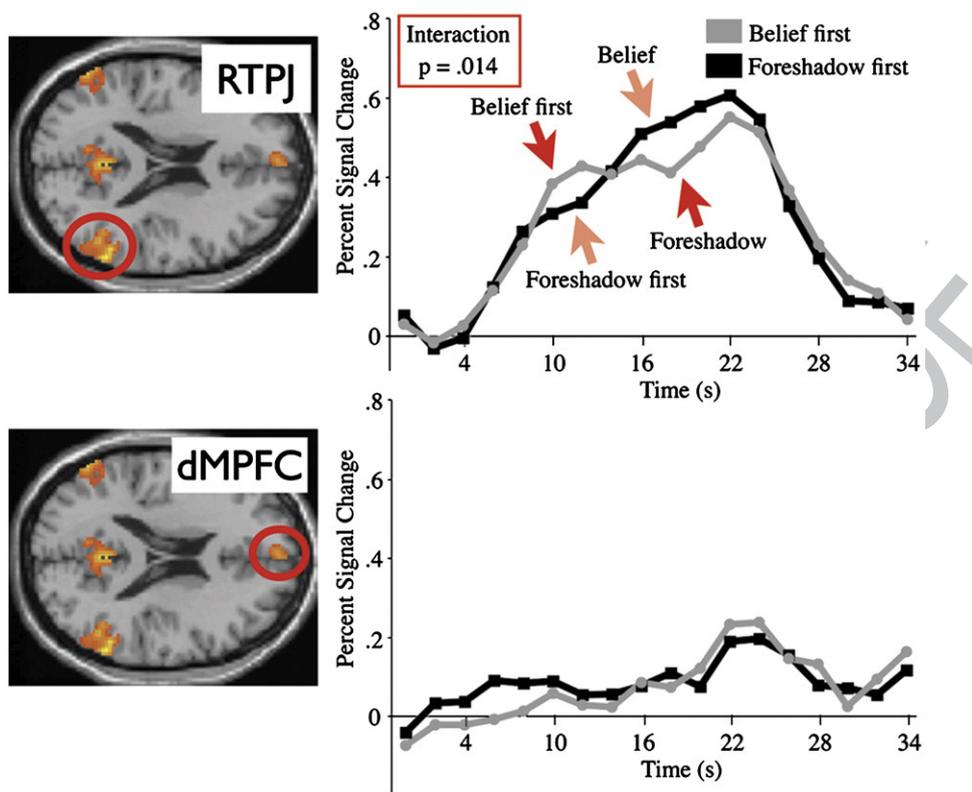


Fig. 2. PSC from rest in the RTPJ (top) and dMPFC (bottom) over time. (Left) Brain regions where the BOLD signal was higher for (nonmoral) stories about beliefs than (nonmoral) stories about physical representations ($N=17$, random effects analysis, $p<0.001$ uncorrected). These data were used to define ROIs: RTPJ (top), dMPFC (bottom). (Right) The PSC in the RTPJ (top) and dMPFC (bottom) during belief-first trials (gray) and foreshadow-first trials (black). Time 1 (10–14 s): belief information was presented during belief-first trials; foreshadow information was presented during foreshadow-first trials. Time 2 (16–20 s): foreshadow information was presented during belief-first trials; belief information was presented during foreshadow-first trials. Time 3 (22–26 s): information was presented about the protagonist’s action and the outcome.

348 ANOVA (Table 3). At time 3, the protagonist’s action, the subject
349 of moral judgment, and its actual outcome were described.

350 (1) RTPJ: Even though no new belief information was presented,
351 the PSC in the RTPJ was significantly above baseline in all eight
352 conditions ($p<0.01$). Also, a significant outcome by belief by order
353 interaction ($F(1,14)=17.0$, $p=0.001$, partial $\eta^2=0.55$) was found,
354 suggesting that the contribution of the factors of belief and outcome
355 depended on the order of information presentation (belief-first vs.
356 foreshadow-first). Each order was therefore analyzed separately.

357 For foreshadow-first, the response showed a main effect of
358 “negative” belief over “neutral” belief ($F(1,14)=9.7$, $p=0.008$) and a
359 belief by outcome interaction ($F(1,14)=11.2$, $p=0.005$; Fig. 3, top
360 panel, as reported in Young et al., 2007, Experiment 2). Planned

comparisons revealed that the PSC was higher for “negative” belief 361
than “neutral” belief in the case of a neutral outcome (“negative”: 0.74, 362
“neutral”: 0.32, $t(14)=4.0$, $p=0.001$), but was not significantly 363
different for “negative” and “neutral” belief in the case of a negative 364
outcome (“negative”: 0.41, “neutral” PSC: 0.51, $t(14)=-1.3$ $p=0.22$). 365
Post-hoc Bonferroni’s t -tests revealed that the PSC for attempted harm 366
was significantly greater than each of the other conditions (unknown 367
harm: $t(14)=2.6$, adjusted $p=0.04$; intentional harm: $t(14)=-3.0$, 368
adjusted $p=0.02$). Consistent with the regions of interest analysis, a 369
random effects whole-brain analysis ($p>0.001$, uncorrected) revealed 370
greater activation for attempted harm (negative belief, neutral 371
outcome) as compared to all-neutral stories in the RTPJ, for this order 372
(average peak voxel coordinates [48 –46 16]). 373

t2.1 Table 2
t2.2 Mean PSC in three ROIs during times 1 and 2 of the moral scenarios

ROI	Mean PSC				Interaction of time×order			
	Belief-first		Foreshadow-first		df	F	p value	Partial η^2
Time 1	Time 2	Time 1	Time 2					
RTPJ	0.41	0.44	0.35	0.54	(1,14)	8.00	0.01	0.36
PC	-0.02	0.06	-0.058	0.16	(1,16)	7.40	0.02	0.32
LTPJ	0.33	0.39	0.35	0.51	(1,15)	5.00	0.04	0.25

t2.5 All three of the regions showed a significant interaction between time (time 1
t2.6 vs. time 2) and order (belief-first vs. foreshadow-first).
t2.7
t2.8

Table 3
Mean PSC in three ROIs during time 3 of the moral scenarios

ROI	Mean PSC (belief, outcome)				Interaction of belief×outcome			
	Neut, Neut	Neut, Neg	Neg, Neut	Neg, Neg	df	F	p value	Partial η^2
RTPJ	0.32	0.51	0.74	0.41				
PC	0.07	0.18	0.29	0.11	(1,16)	7.20	0.02	0.31
LTPJ	0.22	0.4	0.56	0.39	(1,15)	5.00	0.04	0.25

t3.1 All three of the regions showed a significant interaction between negative
t3.2 (Neg) and neutral (Neut) belief and outcome information.
t3.3
t3.4
t3.5
t3.6
t3.7
t3.8

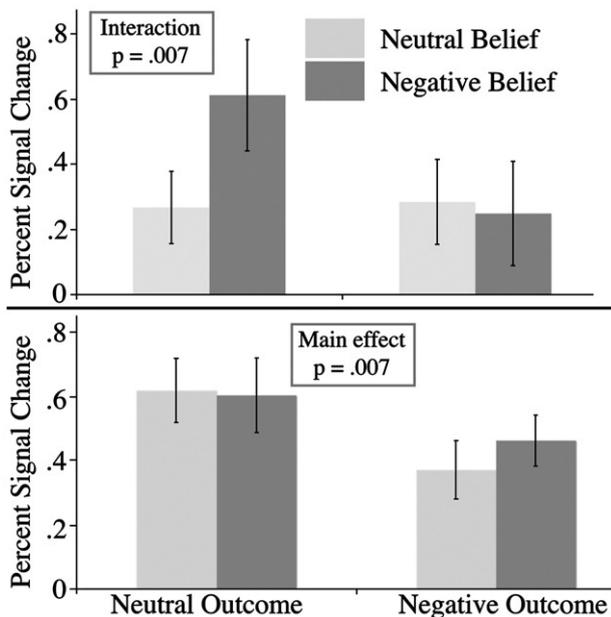


Fig. 3. PSC from rest in the RTPJ during time 3 (22–26 s). Error bars correspond to standard error. (Top) Foreshadow-first trials. (Bottom) Belief-first trials.

A different pattern was observed for belief-first trials (Fig. 3, bottom panel). There was no interaction between belief and outcome. Only a main effect of neutral outcome over negative outcome was significant (neutral: 0.60, negative: 0.36, $F(1,14)=10.2$, $p=0.007$, partial $\eta^2=0.42$). Consistent with this effect, a whole-brain analysis ($p>0.001$, uncorrected) of the overall effect of outcome (neutral over negative) revealed activation in the RTPJ (average peak voxel coordinates [60 –58 20]).

(2) PC and LTPJ: The same $2 \times 2 \times 2$ (outcome [negative vs. neutral] by belief [“negative” vs. “neutral”] by order [belief-first vs. foreshadow-first]) repeated measures ANOVA was conducted for the PC and LTPJ at time 3. A significant belief by outcome by order interaction was found only in the PC ($F(1,16)=5.7$, $p=0.03$). Separate analyses for both orders were performed for both the PC and LTPJ.

A belief by outcome interaction for foreshadow-first was found in the PC ($F(1,16)=7.2$, $p=0.02$). Planned comparisons revealed that the PSC in the PC was higher for “negative” belief than “neutral” belief in the case of a neutral outcome (“negative”: 0.29, “neutral”: 0.07, $t(16)=3.24$, $p=0.005$), but was not significantly different for “negative” and “neutral” belief in the case of a negative outcome (“negative”: 0.11, “neutral” PSC: 0.17, $t(16)=-0.73$, $p=0.48$).

The LTPJ showed the same belief by outcome interaction for foreshadow-first ($F(1,15)=5.0$, $p=0.04$) as well as a main effect of “negative” over “neutral” belief ($F(1,15)=8.5$, $p=0.01$). Planned comparisons revealed that the PSC in the LTPJ was higher for “negative” belief than “neutral” belief in the case of a neutral outcome (“negative”: 0.56, “neutral”: 0.22, $t(15)=3.73$, $p=0.002$), but was not significantly different for “negative” and “neutral” belief in the case of a negative outcome (“negative”: 0.39, “neutral” PSC: 0.40, $t(15)=-0.10$, $p=0.92$).

In contrast to the RTPJ, neither the PC nor the LTPJ showed a main effect of neutral outcome over negative outcome (or any other significant main effects or interactions) for belief-first trials.

(3) MPFC: The belief by outcome by order interaction was not significant in any region of the MPFC. No main effects or interactions

were found for the mMPFC or the vMPFC. The dMPFC, however, 409 showed a main effect of “negative” over “neutral” belief ($F(1,13)=$ 410 9.4, $p=0.01$) for foreshadow-first trials, suggesting a unique role for 411 the dMPFC in processing belief valence for moral judgment. In a 412 previous study using similar stimuli (Young et al., 2007), we observed 413 a similar trend in the dMPFC that did not reach significance 414 (negative>neutral belief, $p<0.1$). However, those results were based 415 on an analysis of only nine individuals. To further investigate the 416 reliability of this effect, we therefore analyzed the response in the 417 dMPFC at the time of integration, across both experiments. A 2×2 418 (belief by outcome) ANOVA ($N=23$) revealed a strong main effect of 419 belief (negative>neutral belief, $F(1,22)=13.3$, $p=0.001$, partial 420 $\eta^2=0.38$), although there was also a significant interaction between 421 belief and outcome ($F(1,22)=11.4$, $p=0.003$, partial $\eta^2=0.34$), 422 similar to that observed in the other regions investigated. 423

General discussion

Moral judgment in the mature state depends on the capacity to 425 attribute beliefs to agents. Both previous and current results suggest 426 that, when belief and outcome information conflict, adult moral 427 judgments are determined primarily by the belief (Cushman, personal 428 communication; Young et al., 2007). Here we distinguish between 429 two cognitive processes associated with belief attribution in moral 430 judgment: the encoding and integration of beliefs. First, belief 431 information is encoded; that is, the relevant belief is detected and 432 represented. Second, belief information is integrated with other 433 relevant information in the construction of moral judgment; belief 434 information is represented in terms of its relation to outcome 435 information. Our results suggest that the same brain regions, the 436 RTPJ, PC, and LTPJ, support both of these belief processes, reflecting 437 a differential response during both encoding and integration phases. 438 The dMPFC, by contrast, appears to process belief valence for moral 439 judgment during the integration phase. Thus, while the RTPJ, PC, 440 and LTPJ are responsible for processing beliefs for moral judgment, 441 the dMPFC is responsible for processing an explicitly morally 442 relevant feature of the action: whether the actor believed he or she 443 was causing harm. Here, we investigate the functional profiles of the 444 response in these regions during moral judgment. 445

The current study reveals neural signatures of the process by 446 which belief information is encoded. This process appears to be 447 supported by the RTPJ, the PC and, to a lesser extent, the LTPJ. 448 Recruitment of these brain regions was observed early in the 449 stimulus, when subjects were first presented with information 450 about the protagonist’s belief. This response was selective for 451 explicit belief information in the current stimulus, as revealed by a 452 significant time by order (belief-first vs. foreshadow-first trials) 453 interaction (cf. Saxe and Wexler, 2005) and consistent with 454 previous research supporting the specific role of these regions but, 455 in particular, the RTPJ, in processing beliefs (Aichorn et al., in 456 press; Fletcher et al., 1995; Gallagher et al., 2000; Gobbin et al., 457 2007; Perner et al., 2006; Saxe and Kanwisher, 2003; Saxe and 458 Powell, 2006; Saxe and Wexler, 2005). Interestingly, the response 459 in these regions at encoding was not influenced by the valence or 460 content of the belief (i.e. “negative” vs. “neutral”) in any of the 461 regions tested. 462

During the integration phase, when subjects were able to make 463 moral judgments of the protagonist’s action and its outcome, the 464 response in the RTPJ, the PC, and LTPJ showed a different functional 465 profile. During this time, no new belief information was added to the 466 story; however, these regions showed above-baseline recruitment 467

468 that differentiated among different moral conditions based on aspects
469 of both the belief and the outcome. The response after the pre-
470 sentation of belief information may reflect the integration of pre-
471 viously presented belief information with other task-relevant
472 information in constructing a coherent moral judgment (Grüneirch,
473 1982; Weiner, 1995; Zelazo et al., 1996). In the context of the current
474 study, outcome information is relevant in two senses: 1) outcome
475 information renders the morally relevant belief true or false, thereby
476 affecting the representation of the belief, and 2) outcome information
477 is independently morally relevant insofar as we judge harms worse
478 than non-harms, and therefore must be reconciled with morally
479 relevant belief information. The current study conflates these two
480 senses of relevance by using outcome information in a moral context,
481 but this distinction should be explored in future studies. Furthermore,
482 while we have focused on the encoding–integration distinction in the
483 context of moral judgment of actions that result in harms or non-
484 harms, it would be of interest to determine what other morally
485 relevant information demands integration with belief information,
486 and whether the same encoding–integration distinction appears in
487 nonmoral domains.

488 The specific functional profile observed during integration dif-
489 fered across brain regions and across stimulus orders. Replicating
490 previous research (Young et al., 2007), we observed a belief by
491 outcome interaction in the RTPJ, PC, and LTPJ when foreshadow
492 information had been presented before belief information. Further-
493 more, the RTPJ response is significantly higher in the case of attempt-
494 ed harm (negative belief, neutral outcome), as compared to each of
495 the other conditions. (Post-hoc comparisons between attempted harm
496 and the other three conditions revealed similar trends in the PC and
497 LTPJ.) A whole-brain random effects group analysis also revealed a
498 greater response uniquely in the RTPJ for attempted harm, contrasted
499 with the all-neutral condition. One interpretation of these results is that
500 moral condemnation depends more heavily on belief information in
501 the absence of a negative outcome. That is, in the case of intentional
502 harm (negative belief, negative outcome), the actor’s causal role in
503 bringing about an actual harm can contribute to moral condemnation.
504 By contrast, in the case of attempted harm (negative belief, neutral
505 outcome), moral condemnation rests solely on the agent’s belief that
506 his or her action will cause harm.

507 However, the response of the RTPJ (though not the PC or the
508 LTPJ) at integration also showed an unexpected interaction with an
509 additional variable: the order of belief and foreshadow information.
510 In contrast to foreshadow-first trials, the RTPJ response at the time
511 of integration of belief-first trials was significantly higher for
512 neutral outcomes than negative outcomes, with no effect of belief
513 valence. Consistent with this main effect, a whole-brain random
514 effects group analysis revealed greater activation in the RTPJ for
515 neutral versus negative outcomes. Moral judgments, by contrast,
516 showed no interaction with order. One explanation for this effect is
517 that participants “double-check” their previously encoded repre-
518 sentation of the protagonist’s beliefs more often, or more deeply,
519 when belief information is presented early, relative to when belief
520 information is presented immediately before the judgment. Future
521 experiments will be necessary to test this hypothesis.

522 It is noteworthy that no aspect of the response in the RTPJ (or the
523 PC or LTPJ) was determined simply by the truth or falsity of the
524 beliefs, as has been suggested by recent work (Sommer et al., 2007).
525 Consistent with our previous study (Young et al., 2007), the current
526 results revealed a significantly above-baseline response in the RTPJ
527 during integration in all eight moral conditions, half of which pre-
528 sented true beliefs; there was no main effect of truth at encoding or

integration. We propose that the moral judgment task of the current 529
study requires reasoning about beliefs, true or false. By contrast, the 530
true belief trials of the “object transfer” task used in the previous 531
research (Sommer et al., 2007) require participants to determine where 532
an observer will look for an object that was “hidden” in full view of the 533
observer. We suggest that this true belief task might not require belief 534
reasoning at all; participants simply have to respond based on the true 535
location of the object (Dennett, 1978). Robust recruitment of the 536
RTPJ, PC, and LTPJ is observed for both true and false beliefs so long 537
as belief reasoning is required by or relevant to the task. 538

We note that our interpretation of the current results is consistent 539
with a specific role for the RTPJ, PC, and LTPJ in belief attribution. 540
Both lesion and imaging studies implicate the RTPJ specifically, 541
however, in another cognitive task: attentional reorienting in 542
response to unexpected stimuli (Corbetta et al., 2000; Mitchell, 543
2007). Nevertheless, the RTPJ response in the current study is best 544
understood as reflecting the processing of belief information, for two 545
reasons. First, attentional reorienting cannot explain the highly 546
selective functional response in the RTPJ. In the encoding phase, for 547
example, belief and outcome information were equally frequent and 548
equally expected, but the RTPJ responded selectively during 549
sentences describing beliefs. Second, a recent study has found that 550
the regions for belief attribution and exogenous attention are 551
neighboring but distinct (Scholz, Triantafyllou, Whitfield-Gabrieli, 552
Brown, Saxe, personal communication). Both individual subject and 553
group analyses revealed less than 8% overlap between the two 554
regions of activation and a reliable separation between the peaks of 555
the two regions: the attention region is located approximately 10 mm 556
superior to the region involved in theory of mind. These results 557
agreed precisely with a recent meta-analysis of 70 published studies 558
that also found that the attention region is 10 mm superior to the 559
region involved in theory of mind (Decety and Lamm, 2007). Given 560
this anatomical separation, the functional localizer approach used in 561
the current study allowed us to identify and then investigate the 562
specific subregion of the RTPJ implicated in theory of mind as well 563
as other regions implicated in theory of mind, i.e. the PC and LTPJ. 564

During both encoding and integration, regions in the MPFC showed 565
a different functional profile. No region in the MPFC distinguished 566
belief from foreshadow information during encoding. Therefore, even 567
though the MPFC is routinely observed in the localizer task, there was 568
no evidence for its specific role in the encoding of beliefs. During the 569
integration phase, however, the dorsal MPFC was selective for the 570
valence of the belief; its response was significantly higher for “nega- 571
tive” than for “neutral” beliefs. In other words, the dMPFC responded 572
more when the protagonist thought that his or her action would cause 573
harm, regardless of whether the action did cause harm. 574

There are two possible accounts for this effect: (1) the dMPFC is 575
responsible for processing belief valence independent of the moral 576
context and (2) the dMPFC is responsible for processing belief 577
valence specifically for moral judgment. We favor the latter account 578
for two reasons. First, belief valence was the dominant factor 579
influencing participants’ moral judgments in this study and other 580
behavioral studies of adult moral judgments (Cushman, personal 581
communication; Young et al., 2007). Second, the dMPFC was 582
recruited differentially for negative over neutral beliefs not during 583
encoding but, rather, only once subjects were able to make moral 584
judgments of the agent’s action, described only during the inte- 585
gration phase. These data suggest a role for the dMPFC in the 586
evaluation of one kind of moral content, specifically, belief valence. 587

These results illuminate prior research suggesting a role for the 588
MPFC in moral judgment (for a review, see Greene and Haidt, 2002; 589

590 Young and Koenigs, in press). Previous research on the neural basis of
 591 moral judgment has focused largely on intentional harm; in all cases the
 592 protagonist knows both that his or her action will cause harm and that
 593 does in fact cause harm by acting. Regions in the MPFC may therefore
 594 have been recruited for representing either actions that produce harmful
 595 *outcomes* or actions performed with harmful *intentions*. The current
 596 results suggest the latter: when subjects are presented with a description
 597 of the critical action, the dMPFC response is sensitive to whether the
 598 actor *thinks* he or she will cause harm by acting.

599 Conclusions

600 The current study reveals the neural basis of at least two distinct
 601 cognitive processes associated directly with theory of mind in moral
 602 judgment, the encoding and the integration of beliefs. Belief encoding
 603 is a stimulus-driven process: the response is based on whether the
 604 current stimulus contains belief information or not. Belief integration
 605 is a relatively stimulus-independent process: prior belief information
 606 is called upon and used in the service of moral judgment. A distinction
 607 between cognitive processes for encoding beliefs versus integrating
 608 beliefs into mature moral judgment is compatible with developmental
 609 research (Baird and Astington, 2004; Baird and Moses, 2001; Zelazo
 610 et al., 1996), and should be further investigated in developmental
 611 cognitive neuroscience. Differential development of function in
 612 theory of mind brain regions, including the RTPJ, PC, and LTPJ, may
 613 coincide with previously reported behavioral changes.

614 Both processes for belief attribution, though, appear to share a neural
 615 substrate in the temporo-parietal junction, bilaterally, and the precuneus.
 616 The medial prefrontal cortex, meanwhile, appears to be uniquely
 617 recruited for processing belief valence, a morally relevant feature of the
 618 action in the context of the task. These results may therefore inform
 619 future research probing the range of contexts both in and beyond the
 620 moral domain that depend on cognitive processes for encoding beliefs,
 621 integrating beliefs, and evaluating the valence of beliefs.

622 Acknowledgments

623 This project was supported by the National Center for Research
 624 Resources (grant P41RR14075), the MIND Institute, and the
 625 Athinoula A. Martinos Center for Biomedical Imaging. R.S. was
 626 supported by MIT and the John Merck Scholars program. L.Y. was
 627 supported by the NSF. Many thanks to Joshua Knobe, Fiery
 628 Cushman, and John Mikhail for comments on an earlier draft of this
 629 manuscript, Jonathan Scholz for technical assistance, and Alexandra
 630 Dickson and Neil Murthy for their help in data collection.

631 Appendix A. Supplementary data

632 Supplementary data associated with this article can be found, in
 633 the online version, at [doi:10.1016/j.neuroimage.2008.01.057](https://doi.org/10.1016/j.neuroimage.2008.01.057).

634 References

- Q1 635 Aichorn, M., Perner, J., Kronblicher, M., Staffen, W., Ladurner, G., in press.
 636 Do visual perspective tasks need theory of mind. *J. Cogn. Neurosci.*
 637 Baird, J.A., Astington, J.W., 2004. The role of mental state understanding in
 638 the development of moral cognition and moral action. *New. Dir. Child.*
 639 *Adolesc. Dev.* 103, 37–49.
 640 Baird, J.A., Moses, L.J., 2001. Do preschoolers appreciate that identical actions
 641 may be motivated by different intentions? *J. Cogn. Dev.* 2, 413–448.
- Borg, J.S., Hynes, C., Van Horn, J., Grafton, S., Sinnott-Armstrong, W., 642
 2006. Consequences, action, and intention as factors in moral 643
 judgments: an fMRI investigation. *J. Cogn. Neurosci.* 18, 803–817. 644
 Ciaramidaro, A., Adenzato, M., Enrici, I., Erk, S., Pia, L., Bara, B.G., 645
 Walter, H., 2007. The intentional network: How the brain reads varieties 646
 of intentions. *Neuropsychologia* 45, 3105–3113. 647
 Corbetta, M., Kincade, J.M., Ollinger, J.M., McAvoy, M.P., Shulman, G.L., 648
 2000. Voluntary orienting is dissociated from target detection in human 649
 posterior parietal cortex. *Nat. Neurosci.* 3, 292–297. 650
 Cushman, F., Young, L., Hauser, M.D., 2006. The role of conscious 651
 reasoning and intuitions in moral judgment: testing three principles of 652
 harm. *Psychol. Sci.* 17, 1082–1089. 653
 Darley, J.M., Zanna, M.P., 1982. Making moral judgment. *Am. Sci.* 70, 654
 515–521. 655
 Decety, J., Lamm, C., 2007. The role of the right temporoparietal junction in 656
 social interaction: how low-level computational processes contribute to 657
 meta-cognition. *Neuroscientist.* 658
 Dennett, D., 1978. Beliefs about beliefs. *Behav. Brain Sci.* 1, 568–570. 659
 Fincham, F.D., Jaspers, J., 1979. Attribution of responsibility to the self and 660
 other in children and adults. *J. Pers. Soc. Psychol.* 37, 1589–1602. 661
 Flavell, J.H., 1999. Cognitive development: children’s knowledge about the 662
 mind. *Ann. Rev. Psychol.* 50, 21–45. 663
 Fletcher, P.C., Happe, F., Frith, U., Baker, S.C., Dolan, R.J., Frackowiak, 664
 R.S.J., Frith, C.D., 1995. Other minds in the brain: a functional imaging 665
 study of “theory of mind” in story comprehension. *Cognition* 57, 109–128. 666
 Gallagher, H.L., Happe, F., Brunswick, N., Fletcher, P.C., Frith, U., Frith, 667
 C.D., 2000. Reading the mind in cartoons and stories: an fMRI study 668
 of ‘theory of mind’ in verbal and nonverbal tasks. *Neuropsychologia* 669
 38, 11–21. 670
 Gobbini, M.I., Koralek, A.C., Bryan, R.E., Montgomery, K.J., Haxby, J.V., 671
 2007. Two takes on the social brain: a comparison of theory of mind 672
 tasks. *J. Cogn. Neurosci.* 19, 1803–1814. 673
 Greene, J.D., Haidt, J., 2002. How (and where) does moral judgment work? 674
Trends Cogn. Sci. 6, 517–523. 675
 Grueneich, R., 1982. The development of children’s integration rules for 676
 making moral judgments. *Child Dev.* 53, 887–894. 677
 Hebble, P.W., 1971. Development of elementary school children’s judgment 678
 of intent. *Child Dev.* 42, 583–588. 679
 Kaniol, R., 1978. Children’s use of intention cues in evaluating behavior. 680
Psychol. Bull. 85, 76–85. 681
 Knobe, J., 2005. Theory of mind and moral cognition: exploring the 682
 connections. *Trends Cogn. Sci.* 9, 357–359. 683
 Mikhail, J., 2007. Universal moral grammar: theory, evidence and the future. 684
Trends Cogn. Sci. 11, 143–152. 685
 Mitchell, J.P., 2007. Activity in right temporo-parietal junction is not 686
 selective for theory-of-mind. *Cereb. Cortex.* 687
 Nelson Le Gall, S.A., 1985. Motive outcome matching and outcome 688
 foreseeability—effects on attribution of intentionality and moral 689
 judgments. *Dev. Psychol.* 21, 332–337. 690
 Nunez, M., Harris, P.L., 1998. Psychological and deontic concepts: separate 691
 domains or intimate connections. *Mind Lang.* 13, 153–170. 692
 Onishi, K., Baillargeon, R., 2005. Do 15-month-old infants understand false 693
 beliefs. *Science* 308, 255–258. 694
 Perner, J., Aichorn, M., Kronblicher, M., Staffen, W., Ladurner, G., 2006. 695
 Thinking of mental and other representations: the roles of left and right 696
 temporo-parietal junction. *Soc. Neurosci.* 1, 245–258. 697
 Piaget, J., 1965/1932. *The Moral Judgment of the Child*. Free Press, New York. 698
 Ruby, P., Decety, J., 2003. What you believe versus what you think they 699
 believe: a neuroimaging study of conceptual perspective-taking. *Eur. J.* 700
Neurosci. 17, 2475–2480. 701
 Saxe, R., Kanwisher, N., 2003. People thinking about thinking people. The 702
 role of the temporo-parietal junction in “theory of mind”. *NeuroImage* 703
 19, 1835–1842. 704
 Saxe, R., Powell, L., 2006. It’s the thought that counts: Specific brain regions 705
 for one component of Theory of Mind. *Psychol. Sci.* 17, 692–699. 706
 Saxe, R., Wexler, A., 2005. Making sense of another mind: The role of the 707
 right temporo-parietal junction. *Neuropsychologia* 43, 1391–1399. 708

- 709 Shultz, T.R., Wright, K., Schleifer, M., 1986. Assignment of moral
710 responsibility and punishment. *Child Dev.* 57, 177–184.
- 711 Siegel, M., Peterson, C.C., 1998. Preschoolers' understanding of lies and
712 innocent and negligent mistakes. *Dev. Psychol.* 34, 332–341.
- 713 Sommer, M., Dohnel, K., Sodian, B., Meinhardt, J., Thoermer, C., Hajak, G.,
714 2007. Neural correlates of true and false belief reasoning. *NeuroImage*
715 35, 1378–1384.
- 716 Vogeley, K., Bussfield, P., Newen, A., Herrmann, S., Happe, F., Falkai, P.,
717 Maier, W., Shaw, N.J., Fink, G.R., Zilles, K., 2001. Mind reading: neural
718 mechanisms of theory of mind and self-perspective. *NeuroImage* 14,
719 170–181.
- 720 Weiner, B., 1995. *Judgments of Responsibility: a Foundation for a Theory of*
721 *Social Conduct*. New York, Guilford Press.
- 722 Wellman, H.M., Cross, D., Watson, J., 2001. Meta-analysis of theory-of-
723 mind development: the truth about false belief. *Child. Dev.* 72, 655–684.
- 724 Wellman, H.M., Larkey, C., Somerville, S.C., 1979. Early development of
725 moral criteria. *Child Dev.* 50, 869–873.
- 726 Young, L., Cushman, F., Hauser, M., Saxe, R., 2007. The neural basis of the
727 interaction between theory of mind and moral judgment. *Proc. Natl.* 727
728 *Acad. Sci.* 104, 8235–8240.
- 729 Young, L., Koenigs, M., in press. Investigating emotion in moral cognition:
730 A review of evidence from functional neuroimaging and neuropsychol-
731 ogy. *British Medical Bulletin*.
- 732 Yuill, N., 1984. Young children's coordination of motive and outcome in
733 judgements of satisfaction and morality. *Br. J. Dev. Psychol.* 2, 73–81.
- 734 Yuill, N., Perner, J., 1988. Intentionality and knowledge in children's
735 judgments of actors responsibility and recipients emotional reaction.
736 *Dev. Psychol.* 24, 358–365.
- 737 Zelazo, P.D., Helwig, C.C., Lau, A., 1996. Intention, act, and outcome in
738 behavioral prediction and moral judgment. *Child Dev.* 67, 2478–2492.

Q4

UNCORRECTED PROOF