



Universiteit  
Leiden  
The Netherlands

## Neural correlates of the motivation to be moral

Nunspeet, F. van

### Citation

Nunspeet, F. van. (2014, May 27). *Neural correlates of the motivation to be moral*. Kurt Lewin Institute Dissertation Series. Ridderprint B.V., Ridderkerk. Retrieved from <https://hdl.handle.net/1887/25829>

Version: Corrected Publisher's Version

License: [Licence agreement concerning inclusion of doctoral thesis in the Institutional Repository of the University of Leiden](#)

Downloaded from: <https://hdl.handle.net/1887/25829>

**Note:** To cite this publication please use the final published version (if applicable).

Cover Page



Universiteit Leiden



The handle <http://hdl.handle.net/1887/25829> holds various files of this Leiden University dissertation.

**Author:** Nunspeet, Félice van

**Title:** Neural correlates of the motivation to be moral

**Issue Date:** 2014-05-27

**Part I**

# **The importance of being moral**



## Chapter 2

# Moral concerns affect implicit prejudice and associated cognitive processes: Behavioral and ERP findings

This chapter is based on:

Van Nunspeet, F., Ellemers, N., Derks, B., & Nieuwenhuis, S. (2014). Moral concerns increase attention and response monitoring during IAT performance: ERP evidence. *Social, Cognitive, and Affective Neuroscience*, 9, 141-149.  
doi:10.1093/scan/nss118



We tend to evaluate people's personal characteristics and behavior along two dimensions: One concerning morality (i.e., how we should behave) and one concerning competence (i.e., how we are able to behave). Behaving according to these dimensions is differentially diagnostic for who we are and how we are perceived: Skowronski and Carlston (1987) showed that for morality negative behaviors are perceived as more diagnostic than positive behaviors, whereas for competence positive behaviors are more diagnostic than negative behaviors. In contrast to behaving incompetently, behaving immorally thus seems to be more indicative of who we are.

Recent research has shown that for people's self-views and the positive evaluation of the group to which they belong moral characteristics are perceived as more important than characteristics concerning competence or sociability (as these are distinct dimensions of social judgment; Leach et al., 2007; in contrast to warmth in which both traits concerning morality and sociability are included; Fiske et al., 2007). Moreover, when people form an impression of a person or a group, they are more interested in information concerning morality traits than traits concerning competence and sociability (Brambilla et al., 2011; Brambilla et al., 2011). Indeed, when people form a first impression within milliseconds, they are more efficient in making inferences about trustworthiness than in making inferences about competence or likeability (Willis & Todorov, 2006).

People seem to be aware that moral judgments are important. For instance, Ellemers et al. (2008) demonstrated that people are inclined to adapt their choice to increase outcomes for the self or for the group to what other group members see as moral than to what other group members see as competent. Moreover, people anticipate being respected by their group members when they adjust their behavior to what the group considers moral (Pagliaro et al., 2011). These findings suggest that morality is of great importance for impression formation and deliberate impression management. We argue that people might also be more inclined to adjust their less deliberate actions (i.e., their implicit behavior) to what is considered moral than to what is considered competent.

In the current study, we examine this prediction using an Implicit Association Test (IAT) that is framed as a test of either individual morality or competence. The

IAT (Greenwald et al., 1998) has been used to measure implicit attitudes towards particular social groups, for example people with dark/white skin or, as in the current study, Muslim/non-Muslim women (see Appendix A). Targets in an IAT consist of stimuli representing social groups that are associated with positive and negative attributes. When people associate stimuli that represent their own (in-)group with positivity and stimuli that represent another (out-)group with negativity, they should respond more quickly and easily to trials that are congruent with these implicit associations than to incongruent combinations (e.g., ingroup stimuli and negative attributes). The IAT assesses the degree to which this is the case, as an indicator of implicit bias.

Whether IAT performance is really implicit and thus uncontrollable is much debated, however. There is much research showing the malleability of implicit attitudes, for example by repeated exposure to admired and disliked individuals (Dasgupta et al., 2009), emotions (Dasgupta & Greenwald, 2001), and several self and social motives (for a review see Blair, 2002). Moreover, it has been shown that the IAT effect is enhanced under stereotype threat (Frantz et al., 2004), but can be diminished when participants have a strategy that helps them to reduce their bias (Fiedler & Bluemke, 2005). In the current research we take advantage of the malleability of the IAT effect: We emphasize the social implications of participants' task performance (i.e., concerning their morality or their competence) and expect that participants to whom the implications concerning morality are emphasized will reduce their negative bias towards Muslim women. More specifically, we hypothesize that these participants will try to inhibit their implicit associations between Muslims and negative attributes, resulting in increased reaction times on congruent trials and thus a smaller IAT effect (consistent with the research of Fiedler & Bluemke [2005]).

Moreover, we are interested in the cognitive processes underlying the motivation to be moral and thus the inhibition of a negative bias on the IAT. Are intentions to behave in line with moral values associated with control of undesirable behavior, or do they influence selective attention that facilitates correct behavior? The current research addresses these questions using event-related brain



potentials (ERPs) associated with perceptual processing and conflict- and error monitoring.

### **Perceptual Attention**

ERPs that are associated with early perceptual processing and more specifically, with selective attention and social categorization are the N1, the P150, and the N2. These components are associated with attention in such a way that increased amplitudes reflect the extent to which attention is directed towards a particular stimulus (e.g., Ito & Urland, 2003). Moreover, research has shown that this attention differs between different social stimuli. For instance, the N1 – a negative deflection occurring around 100 ms after a stimulus is presented – is often larger when viewing stimuli resembling outgroup compared to ingroup members (i.e., black vs. white faces; Ito & Urland, 2003; Kubota & Ito, 2007; however, see Ito & Urland, 2005 for the reversed pattern). The P150, a positive peak that occurs somewhat later (approximately 150-250 ms post-stimulus, therefore also referred to as the P200), is also larger in amplitude for outgroup than for ingroup members (Ito & Urland, 2003; Kubota & Ito, 2007). In contrast to the N1 and P150, the N2 – a negative deflection around 200 ms post-stimulus – is found to be greater for stimuli representing the ingroup compared to the outgroup (Dickter & Bartholow, 2007; Ito & Urland, 2003, 2005). Examination of these components can thus show whether the emphasis on morality attracts greater attention to the group membership of the faces presented in the IAT (which is of importance when the test is said to measure participants' moral values concerning egalitarianism, but not when the test is said to measure competence). Moreover, components related to selective attention and social categorization can also be associated with motivated perception (e.g., the P150/P200; Amodio, 2010). We propose that emphasizing morality increases the motivation to suppress bias towards the outgroup. Although this could lead to diminished social categorization, we hypothesize that social categorization is actually enhanced: People's focus on the different group members should be increased to be sure to respond in line with egalitarian values (i.e., to be able to control implicit bias, as is also seen in research by Amodio, 2010). In other words, we expect to find stronger group-related modulations of the N1, P150, and N2 in the morality condition than in the competence condition.

## Conflict- and Response Monitoring

Because we expect that emphasizing morality motivates people to inhibit their bias, we are also interested in ERPs associated with control. More specifically, conflict- and response monitoring. To assess conflict monitoring, we measure the N450. This is a negative modulation of the ERP signal, typically occurring around 400 ms post-stimulus, when subjects perform incongruent trials. The N450 modulation has been proposed to reflect the occurrence of response conflict (e.g., Nigam et al., 1992; Rebai et al., 1997), and is also evident in incongruent IAT trials (Williams & Themanson, 2011). Because the importance of trial congruency in the IAT may be more evident in the moral than in the competence IAT, and because we expect that control is increased when morality is made salient, we predict that the N450 modulation is larger in the morality compared to the competence condition.

To examine error-monitoring, we assess the error-related negativity (ERN), (Gehring et al., 1993; Nieuwenhuis et al., 2001). The ERN is a negative peak occurring within 100 ms after an erroneous response. The amplitude of the ERN is sensitive to the significance of errors. Hajcak et al. (2005) showed, for example, that ERN amplitude was greater on error trials when fast and accurate responses were associated with a large reward, and when participants' performance was being evaluated by a research assistant. In the current study, we hypothesize that subjects will be more motivated to prevent errors in the morality condition than in the competence condition, because the former might be viewed as a sign of immoral behavior, which is seen as more diagnostic for people's impression formation than incompetent behavior (Skowronski & Carlston, 1987). We therefore predict that erroneous responses will be associated with larger ERN modulations in the morality than in the competence condition.

We conducted two studies to test these predictions. In Study 2.1, we examined our hypothesis that social bias in the IAT is reduced when the test is said to measure participants' morality as opposed to their competence. In Study 2.2, we examined the cognitive processes associated with this reduced bias, as manifested in the ERP components discussed above.

## Study 2.1

### Method

#### Participants.

Sixty-six non-Muslim students from Leiden University (24 males,  $M_{age} = 20.2$  years,  $SD = 1.8$ ) participated in this study for money or course credits.

#### Procedure.

After providing written informed consent, participants performed five blocks of the IAT (Greenwald et al., 1998). Stimuli representing the target concepts consisted of 10 pictures of women without a headscarf (i.e., ingroup pictures) and 10 pictures of women with a headscarf (i.e., outgroup pictures; for details concerning the pretest of these stimuli, see Appendix A). Stimuli that represented positive and negative attributes consisted of 5 pictures of positive scenes, and 5 pictures of negative scenes, selected from the International Affective Picture System (Lang et al., 2005). The stimuli were selected based on the scores for pleasure (i.e., negative pictures with scores  $< 4$  and positive pictures with scores  $> 7$ ).

In a block of congruent trials, ingroup pictures shared the same response key as positive pictures and outgroup pictures the same response key as negative pictures. In a block of incongruent trials, this was the case for ingroup and negative pictures, and outgroup and positive pictures. The order of the congruent and incongruent blocks was counterbalanced across participants. Training blocks (IAT steps 1, 2 and 4) consisted of 26 trials, test blocks (steps 3 and 5) of 156 trials each. Every trial started with a fixation point (with a duration that varied between 500-1500 ms), followed by stimulus presentation (680 ms), and a feedback screen (500 ms). This screen indicated whether participants' response was correct (i.e., green check mark), incorrect (i.e., red cross), or "too late". Participants could not correct their incorrect responses.

***Morality vs. competence task instruction.*** Participants were randomly assigned to an instruction condition. In the morality condition, participants read that the test would indicate their *values* concerning equal treatment of different people. In the competence condition, participants read that the test would indicate how well they are *able* to process new information (for the complete instructions,

see Appendix A). All participants were instructed to respond as quickly and accurately as possible. The test implications were repeated before the start of each test block.

**Checks.** To check that the perceived validity of the IAT did not differ between the conditions, we asked participants after they finished the test to respond to the statement: “My test score can assess what kind of person I am”. Furthermore, two items measured participants’ task engagement: “I think it is important to perform well on this test” and “It does not matter to me what my test score is” [reverse coded], ( $r = .62, p < .001$ ). Participants could respond to each statement on a 7-point scale ranging from “completely disagree” (1) to “completely agree” (7). The experiment took approximately one hour after which participants were debriefed and thanked.

### **The IAT effect.**

The dependent measure was the IAT effect, indicated by the  $D$  score. Based on the scoring algorithm described by Greenwald et al. (2003), this was calculated as the difference in reaction times on incongruent and congruent trials divided by a pooled  $SD$  of all correct trials. We included all trials, replaced error latencies with a replacement value ( $M + 2 SD_{\text{correct}}$ ) and replaced latencies exceeding the maximum response time with the maximum response time of 680 ms.

## **Results and Discussion**

### **Checks.**

As intended, participants in the morality and competence condition did not think differently about the perceived validity of the test;  $M(\text{morality}) = 3.12, SD = 1.65; M(\text{competence}) = 3.24, SD = 1.60; F(1,64) < 1$ . Neither did they differ in their self-reported task engagement:  $M(\text{morality}) = 4.14, SD = 1.00; M(\text{competence}) = 4.24, SD = 1.16; F(1,64) < 1$ .

### **IAT effect.**

Participants showed the standard IAT effect: A negative implicit bias towards the outgroup (i.e., women with a headscarf);  $t(65) = 4.72, p < .001$ . However, this bias was stronger in the competence condition;  $t(32) = 5.40, p < .001$ , than in the morality condition;  $t(32) = 1.77, p = .09$ . More importantly, an ANOVA predicting the  $D$  score from instruction conditions and order of test blocks revealed that the

bias was reduced in the morality, compared to the competence condition;  $M(\text{morality}) = 0.13$ ,  $SD = 0.43$ ;  $M(\text{competence}) = 0.34$ ,  $SD = 0.36$ ;  $F(1,62) = 4.56$ ,  $p = .04$ ,  $\eta^2 = .07$ . The reduced IAT effect was caused by a smaller difference between response times on incongruent and congruent trials in the morality condition: Consistent with previous research (Fiedler & Bluemke, 2005), participants in the morality condition responded somewhat more slowly on congruent correct trials than participants in the competence condition;  $F(1,64) = 3.24$ ,  $p = .08$  (see Figure 2.1)<sup>1</sup>. The percentages of errors did not differ between conditions;  $M(\text{morality}) = 8.81$ ,  $SD = 6.03$ ;  $M(\text{competence}) = 7.73$ ,  $SD = 4.98$ ;  $F(1,64) < 1$ . These behavioral results confirmed our hypothesis that task performance is adjusted when morality is made salient. To test which cognitive processes were modulated to produce the corresponding reduction in IAT score, we conducted Study 2.2.

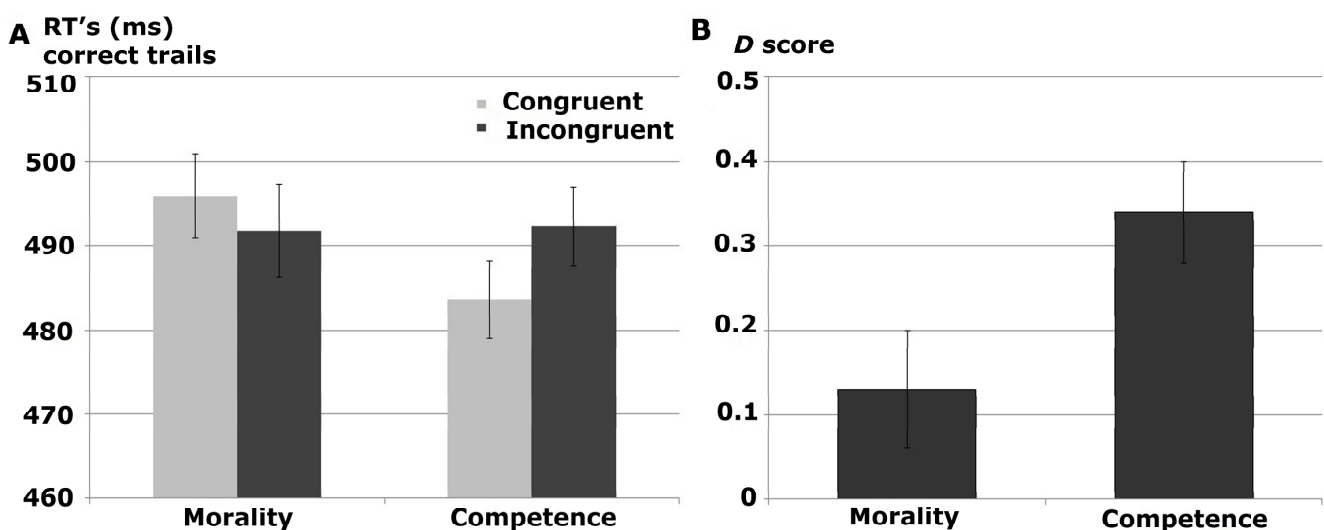


Figure 2.1. Reaction times (in milliseconds) on correct congruent and incongruent trials (A) and the IAT effect, in which error and missed trials are included after they are given a replacement value ( $D$  score; Figure B). Note that the reaction times on incongruent trials are quite fast relative to other IAT studies. This is caused by the limited the presentation time of the stimuli (i.e., participants had to respond within 680ms).

<sup>1</sup> We did not find decreased response times on incongruent trials (which could be expected based on conflict monitoring theory; e.g., Botvinick et al., 2001) because participants had a limited response time.

## Study 2.2

### Method

#### Participants.

Forty-four, healthy, right-handed, non-Muslim students from Leiden University (5 males,  $M_{\text{age}} = 20.4$ ,  $SD = 4.3$ ) provided written informed consent and participated in this study for money or course credits. One participant (morality condition) was excluded from the study due to an outlying IAT score; two participants (morality condition) had to be excluded from EEG analyses because of technical problems. Two more participants (one in each condition) were excluded from statistical analyses of the ERN because they did not make enough errors to reliably quantify this component ( $< 15$ ).

#### Procedure.

Participants performed the IAT as described in Study 2.1, with the following modifications: We inserted a blank screen after the stimulus presentation to ensure that the ERN modulation occurred before the feedback. Each trial thus consisted of a fixation point (500 ms), a stimulus (680 ms), a blank screen (500 ms), and a feedback screen (750 ms). We also increased the number of congruent and incongruent trials from 156 to 300 to enhance the possibility that participants made enough errors to compute a reliable average ERN.

Participants' task engagement was measured with the items from Study 2.1 ( $r = .59$ ,  $p < .001$ ), and we checked whether participants in the morality condition were – as intended – more concerned about the social implications of their performance than participants in the competence condition (i.e., “I am concerned about the impression people might get of me, if they know how I performed on this test”). Moreover, we assessed the internal motivation to respond without prejudice (IMS) scale developed by Plant and Devine (1998; 5 items,  $\alpha = .73$ ; e.g., “I attempt to act in nonprejudiced ways toward women who wear a headscarf because it is personally important to me”; 7-point scale 1 “completely disagree” -7 “completely agree”). Previous research has shown that this internal motivation influences people's ability to regulate biased behavior by conflict-monitoring processes associated with the ERN (Amodio et al., 2008). Thus, to test our prediction that conflict- and error monitoring is enhanced in the morality

compared to the competence condition, we controlled for individual differences in IMS. The total experiment lasted 90 minutes, after which participants were debriefed and thanked.

### **EEG acquisition.**

The EEG was recorded from 19 Ag/AgCl scalp electrodes, and from the left and right mastoids, using a 19-channel Biosemi active-electrode recording system (sampling rate 256 Hz). To assess horizontal and vertical eye movements, electrodes were placed on the outer canthi of the left and right eyes and approximately 1 cm above and below the right eye. EEG activity was recorded using ActiView software, offline data analyses were performed using Brain Vision Analyzer (BVA), and the experiment was controlled by E-prime (v 2.0). The EEG signal was referenced off-line to the average mastoid signal, corrected for ocular and eye-blink artifacts using the method of Gratton et al. (1983), and filtered (1-15 Hz). Single-trial stimulus-locked and response-locked epochs were extracted, ranging from -300 ms to 1000 ms after the event. These epochs were subjected to artifact rejection, then averaged and baseline-corrected by subtracting the average signal value between 200-0 ms pre-stimulus or between 300-50 ms prior to the response. Separate stimulus-locked ERP epochs were created for correct trials with outgroup and ingroup pictures, separately for the congruent and incongruent blocks. Separate response-locked ERP epochs were created for correct and incorrect responses. In an initial analysis, we found no effect of congruency on the ERN. Because participants made few errors on congruent trials, we pooled the congruent and incongruent trials to increase the number of trials averaged for each participant and thus the number of participants included in the ERN analysis.

### **ERP analyses.**

Visual inspection of the data indicated that the N1, P150, N2, and N450 potentials were most evident at the midline electrode sites Fz, FCz, Cz, CPz, and Pz. These ERP components were quantified as the maximum peak amplitude within a time window (N1, 90-110 ms; P150, 100-250 ms; N2, 200-300 ms; N450 325-500 ms). To test the main effects of social categorization and conflict monitoring, we submitted the peak amplitude values to a 5 (electrode site) x 2

(picture type: ingroup/outgroup pictures) x 2 (congruency: congruent/incongruent trials) mixed-model ANOVA.

Visual inspection indicated that the error-related negativity (ERN) was largest at electrodes Fz, FCz, and Cz. To quantify the ERN, we determined the maximal (peak) amplitude of the signal between -50 and 150 ms around the response, separately for correct and incorrect trials. All peak amplitudes were submitted to a 3 (electrode site) x 2 (accuracy: correct/error) mixed-model ANOVA.

Because modulations of the task effects by the instruction manipulation were subtle, subsequent analyses focused on the electrode at which the interaction was most pronounced. The resulting peak-amplitude values were submitted to a mixed-model ANOVA with instruction condition as between-subjects variable and the relevant task factors as within-subject variables. Moreover, to control for individual differences in internal motivation to respond without prejudice, we included IMS score as a covariate in each analysis<sup>2</sup>.

## Results and Discussion

### Checks.

As in Study 2.1, participants in the morality and competence condition did not differ in task engagement;  $M(\text{morality}) = 4.84$ ,  $SD = 0.88$ ;  $M(\text{competence}) = 4.63$ ,  $SD = 0.94$ ;  $F(1,41) < 1$ . Nor did they differ in their internal motivation to respond without prejudice;  $M(\text{morality}) = 4.89$ ,  $SD = 0.82$ ;  $M(\text{competence}) = 5.01$ ,  $SD = 0.66$ ;  $F(1,41) < 1$ . As expected, participants in the morality condition did report to be more concerned about the social implications of their performance than participants in the competence condition;  $M(\text{morality}) = 3.18$ ,  $SD = 1.68$ ;  $M(\text{competence}) = 1.91$ ,  $SD = 1.02$ ;  $F(1,41) = 8.34$ ,  $p = .006$ ,  $\eta^2 = .17$ .

### Behavioral results.

Overall, participants showed the standard IAT effect (i.e., a negative implicit bias towards women with a headscarf);  $t(42) = 5.04$ ,  $p < .001$ . Moreover, this bias was evident in both conditions; morality  $t(20) = 2.52$ ,  $p = .02$ ; competence  $t(21) = 4.68$ ,  $p < .001$ . More importantly, an ANOVA with the  $D$  score based on the first 156 trials in each block as dependent variable, the instruction condition and the

---

<sup>2</sup> Inclusion of the IMS score only changed the results concerning the ERN, as is mentioned in the results section.



order of test blocks as independent variables, and IMS as covariate revealed a difference in the IAT effect between the instruction conditions: As in Study 2.1, the effect was smaller for participants in the morality condition than for participants in the competence condition;  $M(\text{morality}) = 0.13$ ,  $SD = 0.40$ ;  $M(\text{competence}) = 0.42$ ,  $SD = 0.36$ ;  $F(1,39) = 5.86$ ,  $p = .02$ ,  $\eta^2 = .13$ . As can be seen in Figure 2.2, this effect was caused by a smaller difference between response times on incongruent and congruent trials in the morality condition than in the competence condition. More specifically, (and similar to Study 2.1), participants in the morality condition responded somewhat more slowly on congruent trials than participants in the competence condition;  $F(1,41) = 3.06$ ,  $p = .09$ . The percentages of errors did not differ between conditions;  $M(\text{morality}) = 12.36$ ,  $SD = 7.13$ ;  $M(\text{competence}) = 14.25$ ,  $SD = 9.80$ ;  $F(1,41) < 1$ . When we included all trials from each test block (a doubling of trials was needed for computing ERPs), the effect of condition was marginally significant;  $M(\text{morality}) = 0.15$ ,  $SD = 0.27$ ;  $M(\text{competence}) = 0.29$ ,  $SD = 0.29$ ;  $F(1,39) = 3.05$ ,  $p = .09$ . This was caused by a training effect: Participants in both conditions responded faster and made fewer errors on the last 144 trials of each test block, resulting in a similar IAT performance. Although both analyses showed a main effect of the order of test blocks (respectively  $F[1,39] = 23.28$ ,  $p < .001$  and  $F[1,39] = 35.73$ ,  $p < .001$ ), this factor did not interact with instruction condition ( $F$ 's  $< 1$ ).

### **ERP results.**

#### ***Social categorization.***

*N1.* We found the intended main effects of social categorization: The N1 was larger for outgroup pictures ( $M = -5.58 \mu\text{V}$ ,  $S.E. = 0.32$ ) than for ingroup pictures ( $M = -5.26 \mu\text{V}$ ,  $S.E. = 0.30$ );  $F(1,38) = 6.86$ ,  $p = .012$ ,  $\eta^2 = .15$ . Analyses for the FCz electrode confirmed the predicted interaction between instruction condition and picture type;  $F(1,38) = 4.11$ ,  $p = .050$ ,  $\eta^2 = .10$  (see Figure 2.3). The difference between the N1 elicited by outgroup and ingroup pictures was significant in the morality condition ( $F[1,38] = 4.69$ ,  $p = .04$ ,  $\eta^2 = .11$ ), but not in the competence condition ( $F[1,38] < 1$ ).

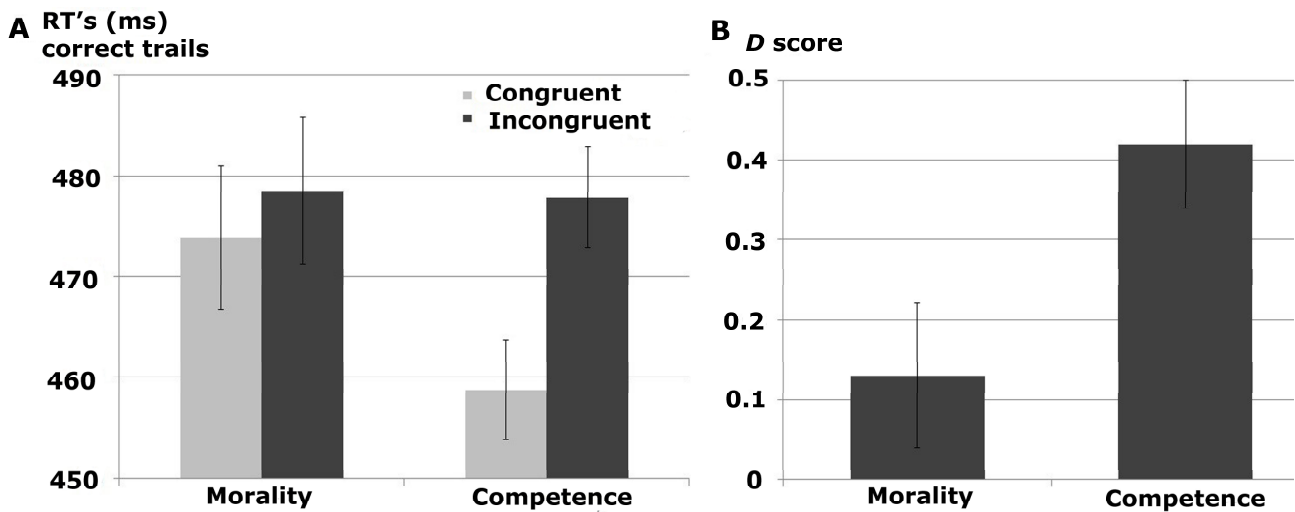


Figure 2.2. Reaction times (in milliseconds) on correct congruent and incongruent trials (A) and the IAT effect in which error and missed trials are included after they are given a replacement value ( $D$  score; Figure B). Note that the reaction times on incongruent trials are quite fast relative to other IAT studies. This is caused by the limited presentation time of the stimuli (i.e., participants had to respond within 680ms).

*P150.* As anticipated, the P150 was larger for outgroup pictures ( $M = 5.22 \mu\text{V}$ ,  $S.E. = 0.52$ ) than for ingroup pictures ( $M = 4.23 \mu\text{V}$ ,  $S.E. = 0.52$ );  $F(1,38) = 39.95$ ,  $p < .001$ ,  $\eta^2 = .51$ . Analyses at Cz showed that, as predicted, there was an interaction effect between instruction condition and picture type;  $F(1,38) = 5.12$ ,  $p = .029$ ,  $\eta^2 = .12$  (see Figure 2.3). The difference in P150 amplitude between outgroup and ingroup pictures was more pronounced in the morality condition ( $F[1,38] = 33.75$ ;  $p < .001$ ,  $\eta^2 = .47$ ), than in the competence condition ( $F[1,38] = 8.51$ ,  $p = .006$ ,  $\eta^2 = .18$ ).

*N2.* The N2 was, as intended, larger for ingroup pictures ( $M = -5.52 \mu\text{V}$ ,  $S.E. = 0.50$ ) than for outgroup pictures ( $M = -4.99 \mu\text{V}$ ,  $S.E. = 0.47$ );  $F(1,38) = 6.93$ ,  $p = .012$ ,  $\eta^2 = .15$ . However, there was no interaction between picture type and instruction condition;  $F(1,38) = 1.08$ ,  $p = .31$ .

### ***Conflict- and error monitoring.***

*N450.* Overall, the N450 was larger for incongruent trials ( $M = -2.22 \mu\text{V}$ ,  $S.E. = 0.39$ ) than for congruent trials ( $M = -1.45 \mu\text{V}$ ,  $S.E. = 0.34$ );  $F(1,38) = 12.51$ ,  $p =$

.001,  $\eta^2 = 0.24$ . Analyses for the CPz electrode confirmed our prediction: Instruction condition interacted with congruency;  $F(1,38) = 4.79$ ,  $p = .035$ ,  $\eta^2 = 0.11$  (see Figure 2.4). The difference in N450 amplitude between incongruent and congruent trials was significant in the morality condition ( $F[1,38] = 16.12$ ,  $p < .001$ ,  $\eta^2 = .30$ ), but not in the competence condition ( $F[1,38] = 1.20$ ,  $p = .28$ ).

*ERN.* As anticipated, the ERN was larger for error trials ( $M = -6.83 \mu\text{V}$ ,  $S.E. = 0.77$ ) than for correct trials ( $M = 1.00 \mu\text{V}$ ,  $S.E. = 0.53$ );  $F(1,36) = 129.08$ ,  $p < .001$ ,  $\eta^2 = 0.78$ . Moreover, accuracy interacted with IMS score;  $F(1,36) = 4.03$ ,  $p = .05$ ,  $\eta^2 = .10$ : A higher internal motivation to respond without prejudice was associated with larger ERN modulations ( $B = -1.46$ ,  $p = .09$ ; see also Amodio et al., 2008). However, more relevant to our current predictions, analyses at Cz showed a marginally significant interaction between accuracy and instruction condition;  $F(1,36) = 3.49$ ,  $p = .070$ ,  $\eta^2 = .09$  (see Figure 2.5)<sup>3</sup>. The difference in ERN amplitude between error and correct trials was somewhat larger in the morality condition ( $M = -11.22 \mu\text{V}$ ,  $S.E. = 1.17$ ;  $F[1,36] = 94.17$ ,  $p < .001$ ,  $\eta^2 = .72$ ) than in the competence condition ( $M = -8.38 \mu\text{V}$ ,  $S.E. = 1.08$ ;  $F[1,36] = 59.74$ ,  $p < .001$ ,  $\eta^2 = .62$ ).

The ERP results are consistent with our expectations that stressing moral implications of the IAT increases social categorization of stimuli and conflict monitoring during the test. More specifically, the emphasis on morality moderates the attention towards outgroup but not ingroup faces (as indexed by increased N1 and P150, but not N2 modulations), and increases the neural response to response conflict and errors in the IAT (as reflected in increased N450 and ERN modulations), suggesting that erroneous responses were perceived as more significant in the morality than in the competence condition.

---

<sup>3</sup> The analysis without IMS as a covariate revealed the same pattern of moderation, but resulted in a non-significant interaction;  $F(1,37) = 2.57$ ,  $p = .12$ . Moreover, as was put forward by an anonymous reviewer, the ERN results were sensitive to changes in the EEG processing settings. For example, shortening the baseline correction period (from 300-50 ms to 200-50 ms prior to the response) reduced the interaction effect between the ERN modulation and instruction;  $F(1,36) = 2.72$ ,  $p = .11$ ,  $\eta^2 = .07$ ; whereas lowering the cutoff score for the high-pass filter (from 1 to 0.1 Hz) made this interaction significant;  $F(1,36) = 4.97$ ,  $p = .03$ ,  $\eta^2 = .12$ .

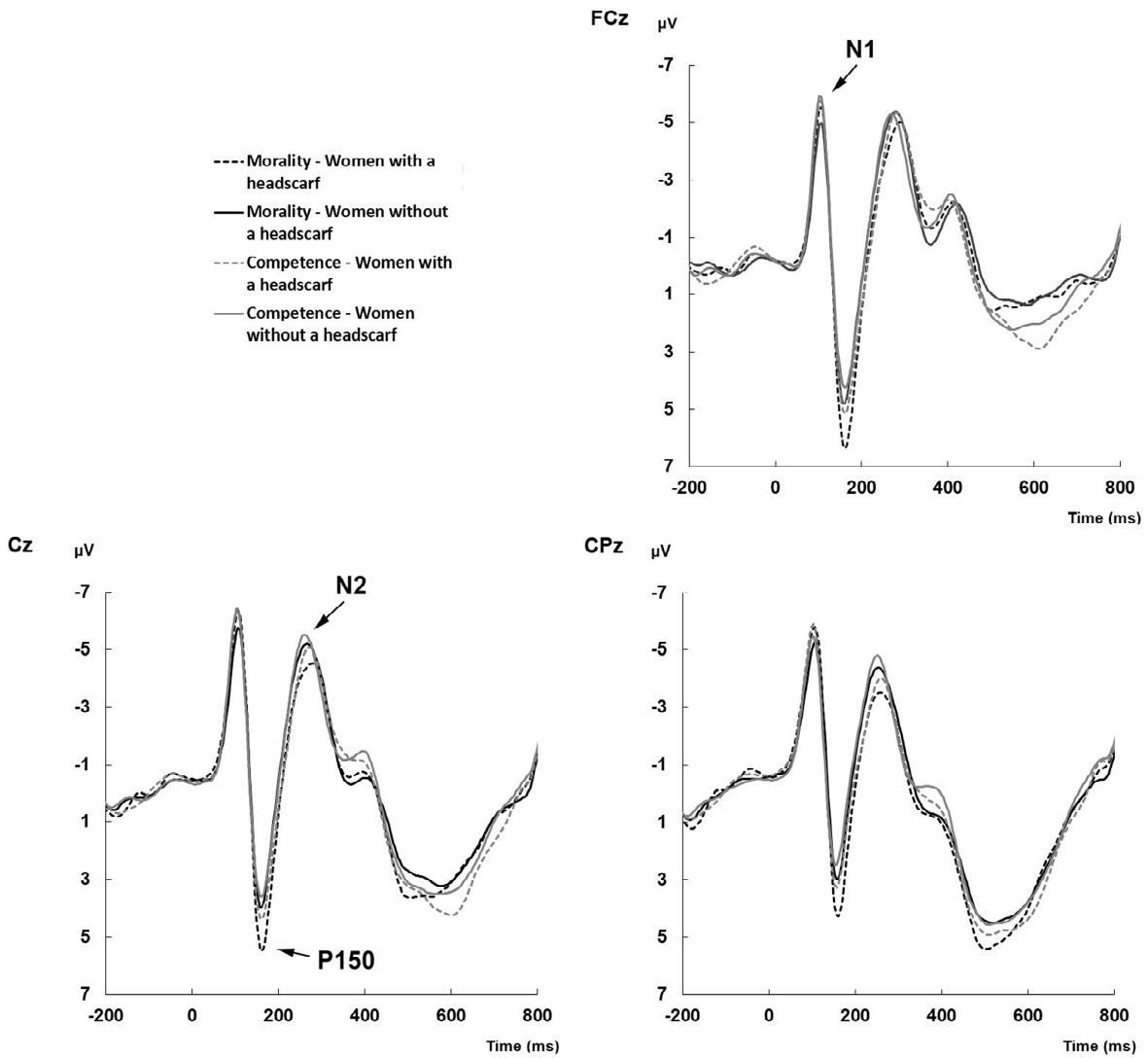


Figure 2.3. The N1, P150 and N2 modulations for pictures of women with and without a headscarf at three central electrodes. The interaction with instruction condition was significant at FCz for the N1, and at Cz for the P150. The interaction did not reach significance for the N2.

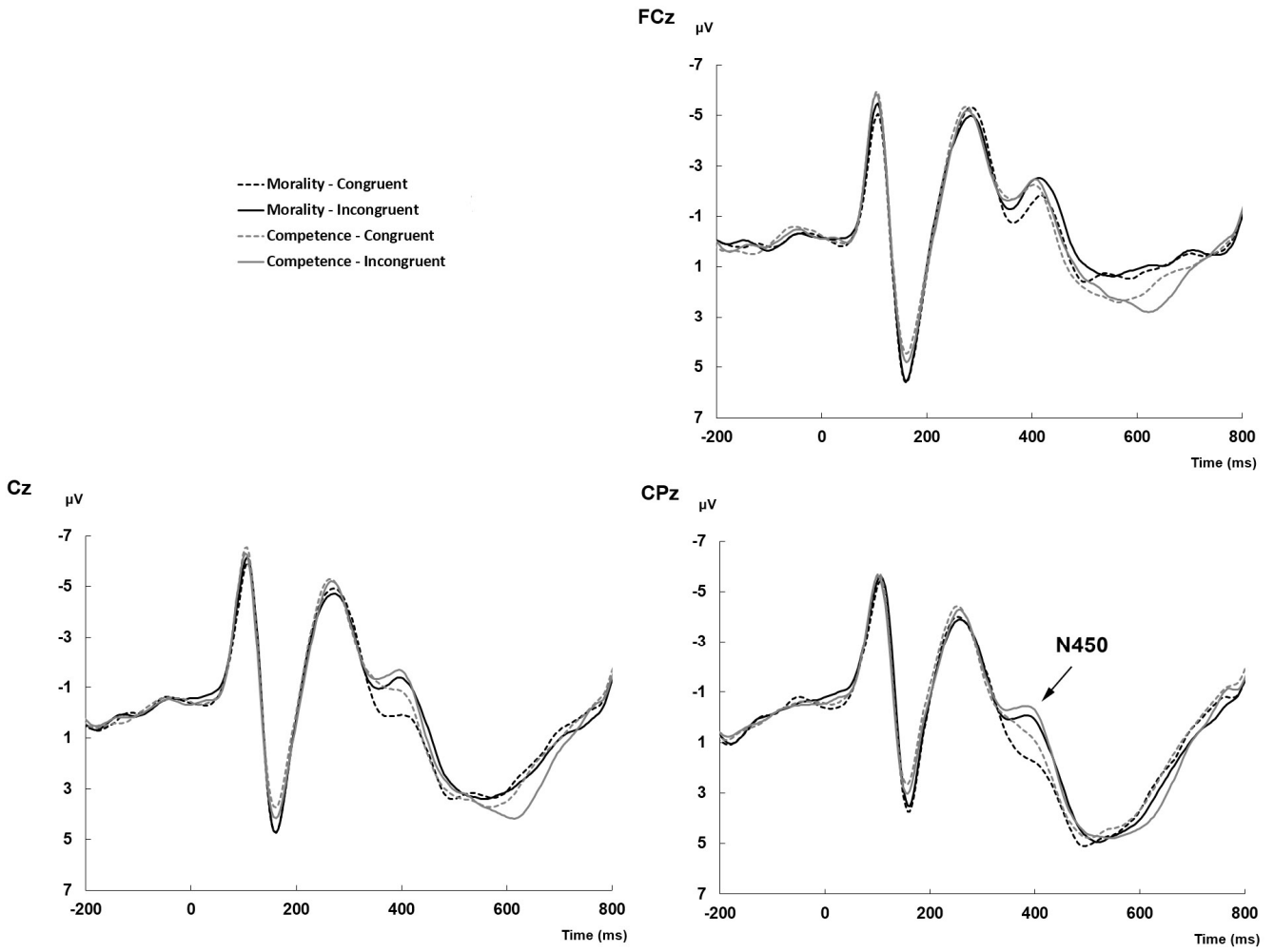


Figure 2.4. The N450 modulations for incongruent and congruent trials at three central electrodes. The interaction with instruction condition was significant at CPz.

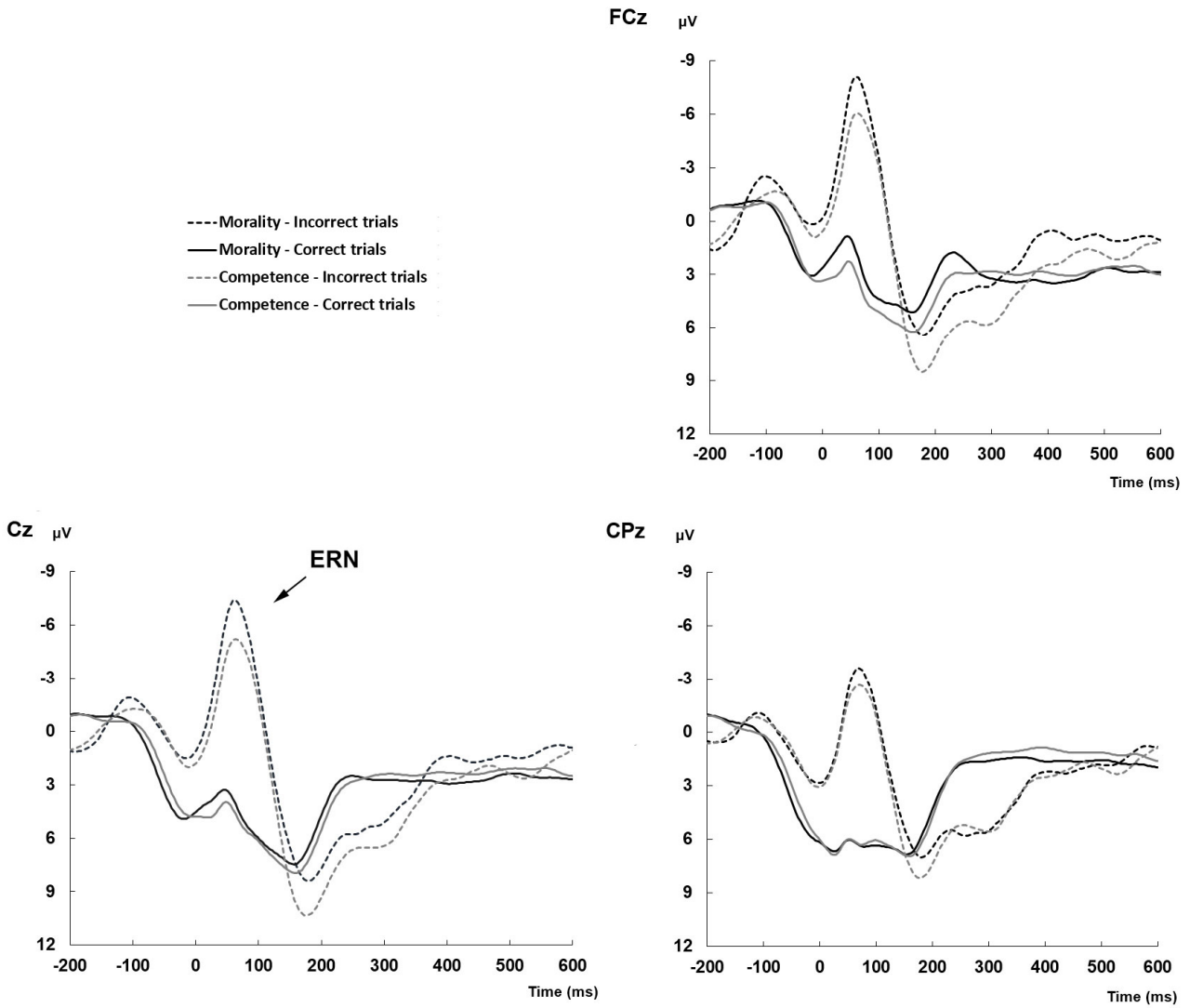


Figure 2.5. The ERN modulations for correct and incorrect trials at three central electrodes. The interaction with instruction condition was marginally significant at Cz.

## General Discussion

Previous research has shown that morality is more important than competence for people's personal and social identity (e.g., Leach et al., 2007), and that morality guides explicit strategic behavior (Ellemers et al., 2008). The present studies extend prior research by showing that morality also impacts on non-explicit aspects of task behavior: People inhibited their negative bias towards Muslim women on an IAT when the test was said to be indicative of their morality (instead of their competence). Our findings thus reveal that participants are able to reduce their implicit bias when given the opportunity to reveal their moral side. This complements prior observations that implicit bias is exacerbated when participants are identified as potential racists (Frantz et al., 2004), and is consistent with research showing that moral appeals induce different physiological and behavioral responses, depending on whether these are framed as ideals or as obligations (Does et al. 2011; 2012).

Importantly, the current research provides insight into the neurobiological mechanisms underlying the differential performance on the moral and competence IAT. Previous research has shown that performance on tasks designed to measure implicit attitudes are associated with (increased) motivated perception (Amodio, 2010) and response monitoring (Amodio, et al., 2008). Additionally, this study reveals that these cognitive processes are activated or enhanced when people's morality is emphasized. More specifically, when morality is emphasized as opposed to competence, people engage in increased social categorization of outgroup faces, and in enhanced conflict- and response monitoring. Because these processes have previously been associated with motivational states (e.g., Amodio, 2010; Hajcak et al., 2005) and because morality has been shown to be more important than competence for impression formation and -management, we interpret these findings as indicating increased motivation of participants in the morality condition to control their bias on the IAT.

The findings concerning increased conflict- and error monitoring during a moral IAT also extend research showing that low levels of implicit bias (often revealed by people with high internal and low external motivation to avoid prejudice) are associated with successful response monitoring (Amodio et al., 2008;

Gonsalkorale et al., 2011). The current results additionally indicate that, regardless of individual differences in internal motivation to respond without prejudice, emphasizing moral values successfully reduces displays of implicit bias. Moreover, our results indicate that emphasizing morality affects not only corrective processes like error monitoring, but affects performance through processes involved in the attention to social stimuli before responses are given.

Although the current research broadens the knowledge of the importance of morality for people's self-identity, we also mentioned that morality is more important than competence for people's *social* identity, and their behavior in groups (Ellemers, et al., 2008; Leach et al., 2007). The question thus remains whether our findings would be affected by for example social evaluation. Further research could address this question by examining whether the emphasis on morality influences people's task performance in the presence of other people and whether this differs between evaluations of ingroup compared to outgroup members.

### **Conclusion**

Our findings extend previous research that demonstrates the importance of morality over competence for people's self-view. In particular, our findings show that people control their implicit responses during a moral task, and reveal how they do that: Emphasizing morality facilitates people's task performance by increasing perceptual attention and conflict- and error monitoring.

### **Acknowledgements**

We thank Ilona Domen, Suzanne Cederhout, Reinier Lagerwerf, Piarella Rodriguez, Lenny van den Beukel, Jelle van Hasselt, and Bart van Wingerde for their help with the data collection, and David Amodio, Guido Band, Stephen Brown, Eveline Crone and Henk van Steenbergen for their advice.