

<sup>1</sup> Good-self based social categorization in perceptual matching

<sup>2</sup> Hu Chuan-Peng<sup>1,2</sup>, Kaiping Peng<sup>3</sup>, & Jie Sui<sup>3,4</sup>

<sup>3</sup> <sup>1</sup> School of Psychology, Nanjing Normal University

<sup>4</sup> <sup>2</sup> Leibniz Institute for Resilience Research, 55131 Mainz, Germany

<sup>5</sup> <sup>3</sup> Tsinghua University, 100084 Beijing, China

<sup>6</sup> <sup>4</sup> University of Aberdeen, Aberdeen, Scotland

<sup>7</sup> Author Note

<sup>8</sup> Hu Chuan-Peng, Leibniz Institute for Resilience Research (LIR). Kaiping Peng,

<sup>9</sup> Department of Psychology, Tsinghua University, 100084 Beijing, China. Jie Sui, School of

<sup>10</sup> Psychology, University of Aberdeen, Aberdeen, Scotland.

<sup>11</sup> Authors contriubtion: HCP, JS, & KP design the study, HCP collected the data,

<sup>12</sup> HCP analyzed the data and drafted the manuscript. KP & JS supported this project.

<sup>13</sup> Correspondence concerning this article should be addressed to Hu Chuan-Peng,

<sup>14</sup> School of Psychology, Nanjing Normal University, Ninghai Road 122, Gulou District,

<sup>15</sup> 210024 Nanjing, China. E-mail: hcp4715@gmail.com

16

## Abstract

17 To navigate in a complex social world, individuals are constantly evaluating others' moral  
18 character. Also, they are managing a moral self-view that is aligned with their goals.

19 Though moral character in person perception and moral self-enhancement had been  
20 extensively studied, the perceptual process of moral character is unkown, we examined  
21 the influence of immediately acquired moral information on perceptual matching processing  
22 by using social associative learning paradigm (self-tagging paradigm). In a series of  
23 experiments, participants learned the concept of moral character and visual cues (shapes)  
24 and then perform a perceptual matching task. The results showed that when geometric  
25 shapes, without soical meaning, that associated with good moral character were prioritized.

26 This patterns of results were robust when we change different semantic words or using  
27 behavioral history as an proxy of mroal character. Also, this patterns were robust in both  
28 simutanously presentation and sequential presentation. We then examined two competing  
29 explanation for this effect: value-based prioritization or social-categorization based  
30 prioritization. We manipulated the identity of different moral character explicitly and  
31 found that the good moral character effect was strong when for the self-referential  
32 conditions but not for other-referential condition. We further tested the good-self based  
33 social categorization by presenting the identity or moral character information as  
34 task-irrelevant stimuli, so that we can distinguish between the unique good-self hypothesis  
35 and a more general good-person based social categorization hypothesis. The evidence  
36 suggested that human are more likely has a good-person based categorization instead of a  
37 unique good-self. Finally, we explored whether the positivity effect only exist in moral  
38 domain and found that this effect was not limited to moral domain but also aesthetic  
39 domain, but not affective valence *per se*. Exploratory analyses on task-questionnaire  
40 relationship found that there are weak correlation between self-bad distance and behavioral  
41 pattern. These results suggest that there exist a social categorization in perceptual  
42 decision-making, which is based on personal traits (moral character) but not affective

<sup>43</sup> valence.

<sup>44</sup> *Keywords:* Perceptual decision-making, Self positivity bias, morality

<sup>45</sup> Word count: X

<sup>46</sup> Good-self based social categorization in perceptual matching

<sup>47</sup> **Introduction**

<sup>48</sup> [sentences in bracket are key ideas]

<sup>49</sup> [Morality is the central of human social life]. People experience a substantial amount  
<sup>50</sup> of moral events in everyday life (e.g., Hofmann, Wisneski, Brandt, & Skitka, 2014). When  
<sup>51</sup> experiencing these events, it always involves judging “good” or “bad.” Judging “good”  
<sup>52</sup> vs. “bad” also appeared implicitly in judging “right” or “wrong,” i.e., moral character  
<sup>53</sup> (Uhlmann, Pizarro, & Diermeier, 2015). Similarly, moral character is a basic dimension of  
<sup>54</sup> person perception (Abele, Ellemers, Fiske, Koch, & Yzerbyt, 2020; Goodwin, 2015;  
<sup>55</sup> Goodwin, Piazza, & Rozin, 2014; Willis & Todorov, 2006) and the most important aspect  
<sup>56</sup> to evaluate the continuity of identity (Strohminger, Knobe, & Newman, 2017).

<sup>57</sup> Given the importance of moral character, to successfully navigate in a social world, a  
<sup>58</sup> person needs to both accurately evaluate others’ moral character and behave in a way that  
<sup>59</sup> she/he is perceived as a moral person, or at least not a morally bad person. Maintaining a  
<sup>60</sup> moral self-view is as important as making judgment about others’ moral character  
<sup>61</sup> (Ellemers, Toorn, Paunov, & Leeuwen, 2019). Moral character is studied extensively both  
<sup>62</sup> in person perception (Abele, Ellemers, Fiske, Koch, & Yzerbyt, 2020; Goodwin, 2015;  
<sup>63</sup> Goodwin, Piazza, & Rozin, 2014; Willis & Todorov, 2006) and moral self-view (Klein &  
<sup>64</sup> Epley, 2016; Monin & Jordan, 2009; Strohminger, Knobe, & Newman, 2017; Tappin &  
<sup>65</sup> McKay, 2017). Recent theorists are trying to bring them together and emphasize a  
<sup>66</sup> person-centered moral psychology(Uhlmann, Pizarro, & Diermeier, 2015). In this new  
<sup>67</sup> perspective, role of perceives’ self-relevance in morality has also been studied (e.g., Waytz,  
<sup>68</sup> Dungan, & Young, 2013).

<sup>69</sup> To date, however, as Freeman and Ambady (2011) put it, studies in the perception of  
<sup>70</sup> moral character didn’t try to explain the perceptual process, rather, they are trying to

71 explain the higher-order social cognitive processes that come after. Essentially, these  
72 studies are perception of moral character without perceptual process. Without knowledge  
73 of perceptual processes, we can not have a full picture of how moral character is processed  
74 in our cognition. As an increasing attention is paid to perceptual process underlying social  
75 cognition, it's clear that perceptual processes are strongly influenced by social factors, such  
76 as group-categorization, stereotype (Stolier & Freeman, 2016; see Xiao, Coppin, & Bavel,  
77 2016). Given the importance of moral character and that moral character related  
78 information has strong influence on learning and memory (Carlson, Maréchal, Oud, Fehr,  
79 & Crockett, 2020; Stanley & De Brigard, 2019), one might expect that moral character  
80 related information could also play a role in perceptual process.

81 To explore the perceptual process of moral character and the underlying mechanism,  
82 we conducted a series of experiments to explore (1) whether we can detect the influence of  
83 moral character information on perceptual decision-making in a reliable way, and (2)  
84 potential explanations for the effect. In the first four experiment, we found a robust effect  
85 of good-person prioritization in perceptual decision-making. The we explore the potential  
86 explanations and tested value-based prioritization versus self-relevance-based prioritization  
87 (social-categorization (Turner, Hogg, Oakes, Reicher, & Wetherell, 1987; Turner, Oakes,  
88 Haslam, & McGarty, 1994)). These results suggested that people may categorize self and  
89 other based on moral character; in these categorizations, the core self, i.e., the good-self, is  
90 always prioritized.

## 91 Perceptual process of moral character

92 [exp1a, b, c, and exp2]

93 [using associative learning task to study the moral character's influence on  
94 perception] Though it is theoretically possible that moral character related information  
95 may be prioritized in perceptual process, no empirical studies had directly explored this

96 possibility. There were only a few studies about the temporal dynamics of judging the  
97 trustworthiness of face (e.g., Dzhelyova, Perrett, & Jentzsch, 2012), but trustworthy is not  
98 equal to morality.

99 One difficulty of studying the perceptual process of moral character is that moral  
100 character is an inferred trait instead of observable feature. usually, one needs necessary  
101 more sensory input, e.g., behavior history, to infer moral character of a person. For  
102 example, Anderson, Siegel, Bliss-Moreau, and Barrett (2011) asked participant to first  
103 study the behavioral description of faces and then asked them to perform a perceptual  
104 detection task. They assumed that by learning the behavioral description of a person  
105 (represented by a face), participants can acquire the moral related information about faces,  
106 and the associations could then bias the perceptual processing of the faces (but see Stein,  
107 Grubb, Bertrand, Suh, and Verosky (2017)). One drawback of this approach is that  
108 participants may differ greatly when inferring the moral character of the person from  
109 behavioral descriptions, given that notion what is morality itself is varying across  
110 population Jones et al. (2020) and those descriptions and faces may themselves are  
111 idiosyncratic, therefore, introduced large variation in experimental design.

112 An alternative is to use abstract semantic concepts. Abstract concepts of moral  
113 character are used to describe and represent moral characters. These abstract concepts  
114 may be part of a dynamic network in which sensory cue, concrete behaviors and other  
115 information can activate/inhibit each other (e.g., aggressiveness) (Amodio, 2019; Freeman  
116 & Ambady, 2011). If a concept of moral character (e.g., good person) is activated, it  
117 should be able to influence on the perceptual process of the visual cues through the  
118 dynamic network, especially when the perceptual decision-making is about the concept-cue  
119 association. In this case, abstract concepts of moral character may serve as signal of moral  
120 reputation (for others) or moral self-concept. Indeed, previous studies used the moral  
121 words and found that moral related information can be perceived faster Firestone & Scholl  
122 (2015). If moral character is an important in person perception, then, just as those other

<sup>123</sup> information such as races and stereotype (see Xiao, Coppin, & Bavel, 2016), moral  
<sup>124</sup> character related concept might change the perceptual processes.

<sup>125</sup> To investigate the above possibility, we used an associative learning paradigm to  
<sup>126</sup> study how moral character concept change perceptual decision-making. In this paradigm,  
<sup>127</sup> simple geometric shapes were paired with different words whose dominant meaning is  
<sup>128</sup> describing the moral character of a person. Participants first learn the associations between  
<sup>129</sup> shapes and words, e.g., triangle is a good-person. After building direct association between  
<sup>130</sup> the abstract moral characters and visual cues, participants then perform a matching task  
<sup>131</sup> to judge whether the shape-word pair presented on the screen match the association they  
<sup>132</sup> learned. This paradigm has been used in studying the perceptual process of self-concept,  
<sup>133</sup> but had also proven useful in studying other concepts like social group (F. E. Enock,  
<sup>134</sup> Hewstone, Lockwood, & Sui, 2020; F. Enock, Sui, Hewstone, & Humphreys, 2018). By  
<sup>135</sup> using simple and morally neutral shapes, we controlled the variations caused by visual cues.

<sup>136</sup> Our first question is, whether the words used the in the associative paradigm is really  
<sup>137</sup> related to the moral character? As we reviewed above, previous theories, especially the  
<sup>138</sup> interactive dynamic theory, would support this assumption. To validate that moral  
<sup>139</sup> character concepts activated moral character as a social cue, we used four experiments to  
<sup>140</sup> explore and validate the paradigm. The first experiment directly adopted associative  
<sup>141</sup> paradigm and change the words from “self,” “friend,” and “stranger” to “good-person,”  
<sup>142</sup> “neutral-person,” and “bad-person.” Then, we change the words to the ones that have  
<sup>143</sup> more explicit moral meaning (“kind-person,” “neutral-person,” and “evil-person”). Then,  
<sup>144</sup> as in Anderson, Siegel, Bliss-Moreau, and Barrett (2011), we asked participant to learn the  
<sup>145</sup> association between three different behavioral histories and three different names, and then  
<sup>146</sup> use the names, as moral character words, for associative learning. Finally, we also tested  
<sup>147</sup> that simultaneously present shape-word pair and sequentially present word and shape  
<sup>148</sup> didn’t change the pattern. All of these four experiments showed a robust effect of moral  
<sup>149</sup> character, that is, the positive moral character associated stimuli were prioritized.

150 **Morality as a social-categorization?**

151 [possible explanations: person-based self-categorization vs. stimuli-based valence] The  
152 robust pattern from our first four experiment suggested that there are some reliable  
153 mechanisms underneath the effect. One possible explanation is the value-based attention,  
154 which suggested that valuable stimuli is prioritized in our low-level cognitive processes.  
155 Because positive moral character is potentially rewarding, e.g., potential cooperators, it is  
156 valuable to individuals and therefore being prioritized. There are also evidence consistent  
157 with this idea []. For example, XXX found that trustworthy faces attracted attention more  
158 than untrustworthy faces, probably because trustworthy faces are more likely to be the  
159 collaborative partners subsequent tasks, which will bring reward. This explanation has an  
160 implicit assumption, that is, participants were automatically viewing these stimuli as  
161 self-relevant (Juechems & Summerfield, 2019; Reicher & Hopkins, 2016) and  
162 threatening/rewarding because of their semantic meaning. In this explanation, we will view  
163 the moral concept, and the moral character represented by the concept, as objects and only  
164 judge whether they are rewarding/threatening or potentially rewarding/threatening to us.

165 Another possibility is that we will perceive those moral character as person and  
166 automatic categorize whether they are ingroup or ougroup, that is, the social  
167 categorization process. This account assumed that moral character served as a way to  
168 categorize other. In the first four experiments' situation, the identity of the moral  
169 character is ambiguous, participants may automatically categorize morally good people as  
170 ingroup and therefore preferentially processed these information.

171 However, the above four experiments can not distinguish between these two  
172 possibilities, because the concept “good-person” can both be rewarding and be categorized  
173 as ingroup member, and previous studies using associative learning paradigm revealed that  
174 both rewarding stimuli (e.g., Sui, He, & Humphreys, 2012) and in-group information [F.  
175 Enock, Sui, Hewstone, and Humphreys (2018); enock\_overlap\_2020] are prioritized.

[Distinguish two explanations by make self salient, exp3a, 3b, 6b] Though these two frameworks can both account for the positivity effect found in first four experiments (i.e., prioritization of “good-person,” but not “neutral person” and “bad person”), they have different prediction if the experiment design include both identity and moral valence where the valence (good, bad, and neutral) conditions can describe both self and other. In this case the identity become salient and participants are less likely to spontaneously identify a good-person other than self as the extension of self, but the value of good-person still exists. Actually, the rewarding value of good-other might be even stronger than good-self because the former indicate potential cooperation and material rewards, but the latter is more linked to personal belief. This means that the social categorization theory predicts participants prioritize good-self but not good-other, while reward-based attention theory predicts participants are both prioritized. Also, as in Hu, Lan, Macrae, and Sui (2020), people may also only identify with good-self instead of bad self. That is, people will show a unique pattern of self-identification: only good-self is identified as “self” while all the others categories were excluded.

In exp 3a, 3b, and 6b, we found that (1) good-self is always faster than neutral-self and bad-self, but good-other only have weak to null advantage to neutral-other and bad-other. which mean the social categorization is self-centered. (2) good-self’s advantage over good other only occur when self- and other- were in the same task. i.e. the relative advantage is competition based instead of absolute. These three experiments suggest that people more like to view the moral character stimuli as person and categorize good-self as an unique category against all others. A mini-meta-analysis showed that there was no effect of valence when the identity is other. This results showed that value-based attention is not likely explained the pattern we observed in first four experiments. Why good-self is prioritized is less clear. Besides the social-categorization explanation, it’s also possible that good self is so unique that it is prioritized in all possible situation and therefore is not social categorization per se.

[what we care? valence of the self exp4a or identity of the good exp4b?] We go further to disentangle the good-self complex: is it because the special role of good-self or because of social categorization. We designed two complementary experiments. in experiment 4a, participants only learned the association between self and other, the words “good-person,” “neutral person,” and “bad person” were presented as task-irrelevant stimuli, while in experiment 4b, participants learned the associations between “good-person,” “neutral-person,” and “bad-person,” and the “self” and “other” were presented as task-irrelevant stimuli. These two experiment can be used to distinguish the “good-self” as anchor account and the “good-self-based social categorization” account. If good-self as an anchor is true, then, in both experiment, good-self will show advantage over other stimuli. More specifically, in experiment 4a, in the self condition, there will be advantage for good as task-irrelevant condition than the other two self conditions; in experiment 4b, in the good condition, there will be an advantage for self as task-irrelevant condition over other as task-irrelevant condition. If good-self-based social categorization if true, then, the prioritization effect will depends on whether the stimuli can be categorized as the same group of good-self. More specifically, in experiment 4a, there will be good-as-task-irrelevant stimuli than other condition in self conditions, this prediction is the same as the “good-self as anchor” account; however, for experiment 4b, there will be no self-as-task-irrelevant stimuli than other-as-task-irrelevant condition.

[Good self in self-reported data] As an exploration, we also collected participants' self-reported psychological distance between self and good-person, bad-person, and neutral-person, moral identity, moral self-image, and self-esteem. All these data are available (see Liu et al., 2020). We explored the correlation between self-reported distance and these questionnaires as well as the questionnaires and behavioral data. However, given that the correlation between self-reported score and behavioral data has low correlation (Dang, King, & Inzlicht, 2020), we didn't expect a high correlation between these self-reported measures and the behavioral data.

[whether categorize self as positive is not limited to morality] Finally, we explored the pattern is generalized to all positive traits or only to morality. We found that self-categorization is not limited to morality, but a special case of categorization in perpetual processing.

Key concepts and discussing points:

**Self-categories** are cognitive groupings of self and some class of stimuli as identical or different from some other class. [Turner et al.]

**Personal identity** refers to self-categories that define the individual as a unique person in terms of his or her individual differences from other (in-group) persons.

**Social identity** refers to the shared social categorical self (“us” vs. “them”).

**Variable self:** Who we are, how we see ourselves, how we define our relations to others (indeed whether they are construed as ‘other’ or as part of the extended ‘we’ self) is different in different settings.

**Identification:** the degree to which an individual feels connected to an ingroup or includes the ingroup in his or her self-concept. (self is not bad; )

Morality as a way for social-categorization (McHugh, McGann, Igou, & Kinsella, 2019)? People are more likely to identify themselves with trustworthy faces (Verosky & Todorov, 2010) (trustworthy faces has longer RTs).

What is the relation between morally good and self in a semantic network (attractor network) (Freeman & Ambady, 2011).

How to deal with the *variable self* (self-categorization theory) vs. *core/true/authentic self* vs. *self-enhancement*

**Limitations:** The perceptual decision-making will show certain pattern under certain task demand. In our case, it's the forced, speed, two-option choice task.

254

## Disclosures

255 We reported all the measurements, analyses, and results in all the experiments in the  
256 current study. Participants whose overall accuracy lower than 60% were excluded from  
257 analysis. Also, the accurate responses with less than 200ms reaction times were excluded  
258 from the analysis.

259 All the experiments reported were not pre-registered. Most experiments (1a ~ 6b,  
260 except experiment 3b) reported in the current study were first finished between 2014 to  
261 2016 in Tsinghua University, Beijing, China. Participants in these experiments were  
262 recruited in the local community. To increase the sample size of experiments to 50 or more  
263 (Simmons, Nelson, & Simonsohn, 2013), we recruited additional participants in Wenzhou  
264 University, Wenzhou, China in 2017 for experiment 1a, 1b, 4a, and 4b. Experiment 3b was  
265 finished in Wenzhou University in 2017. To have a better estimation of the effect size, we  
266 included the data from two experiments (experiment 7a, 7b) that were reported in Hu,  
267 Lan, Macrae, and Sui (2020) (See Table S1 for overview of these experiments).

268 All participant received informed consent and compensated for their time. These  
269 experiments were approved by the ethic board in the Department of Tsinghua University.

270

## General methods

### 271 Design and Procedure

272 This series of experiments studied the perceptual process of moral character, using  
273 the social associative learning paradigm (or tagging paradigm)(Sui, He, & Humphreys,  
274 2012), in which participants first learned the associations between geometric shapes and  
275 labels of person with different moral character (e.g., in first three studies, the triangle,  
276 square, and circle and good person, neutral person, and bad person, respectively). The  
277 associations of the shapes and label were counterbalanced across participants. After

remembered the associations, participants finished a practice phase to familiar with the task, in which they viewed one of the shapes upon the fixation while one of the labels below the fixation and judged whether the shape and the label matched the association they learned. When participants reached 60% or higher accuracy at the end of the practicing session, they started the experimental task which was the same as in the practice phase.

The experiment 1a, 1b, 1c, 2, and 6a shared a 2 (matching: match vs. nonmatch) by 3 (moral character: good person vs. neutral person vs. bad person) within-subject design. Experiment 1a was the first one of the whole series studies and found the prioritization of stimuli associated with good-person. To confirm that it is the moral character that caused the effect, we further conducted experiment 1b, 1c, and 2. More specifically, experiment 1b used different Chinese words as label to test whether the effect only occurred with certain familiar words. Experiment 1c manipulated the moral valence indirectly: participants first learned to associate different moral behaviors with different neutral names, after remembered the association, they then performed the perceptual matching task by associating names with different shapes. Experiment 2 further tested whether the way we presented the stimuli influence the effect of valence, by sequentially presenting labels and shapes. Note that part of participants of experiment 2 were from experiment 1a because we originally planned a cross task comparison. Experiment 6a, which shared the same design as experiment 2, was an EEG experiment which aimed at exploring the neural correlates of the effect. But we will focus on the behavioral results of experiment 6a in the current manuscript.

For experiment 3a, 3b, 4a, 4b, 6b, 7a, and 7b, we included self-reference as another within-subject variable in the experimental design. For example, the experiment 3a directly extend the design of experiment 1a into a 2 (matchness: match vs. nonmatch) by 2 (reference: self vs. other) by 3 (moral valence: good vs. neutral vs. bad) within-subject design. Thus in experiment 3a, there were six conditions (good-self, neutral-self, bad-self, good-other, neutral-other, and bad-other) and six shapes (triangle, square, circle, diamond,

305 pentagon, and trapezoids). The experiment 6b was an EEG experiment extended from  
306 experiment 3a but presented the label and shape sequentially. Because of the relatively  
307 high working memory load (six label-shape pairs), experiment 6b were conducted in two  
308 days: the first day participants finished perceptual matching task as a practice, and the  
309 second day, they finished the task again while the EEG signals were recorded. Experiment  
310 3b was designed to separate the self-referential trials and other-referential trials. That is,  
311 participants finished two different types of block: in the self-referential blocks, they only  
312 responded to good-self, neutral-self, and bad-self, with half match trials and half  
313 non-match trials; in the other-reference blocks, they only responded to good-other,  
314 neutral-other, and bad-other. Experiment 7a and 7b were designed to test the cross task  
315 robustness of the effect we observed in the aforementioned experiments (see, Hu, Lan,  
316 Macrae, & Sui, 2020). The matching task in these two experiments shared the same design  
317 with experiment 3a, but only with two moral character, i.e., good vs. bad. We didn't  
318 include the neutral condition in experiment 7a and 7b because we found that the neutral  
319 and bad conditions constantly showed non-significant results in experiment 1 ~ 6.

320       Experiment 4a and 4b were design to explore the mechanism behind the  
321 prioritization of good-self. In 4a, we used only two labels (self vs. other) and two shapes  
322 (circle, square). To manipulate the moral valence, we added the moral-related words within  
323 the shape and instructed participants to ignore the words in the shape during the task. In  
324 4b, we reversed the role of self-reference and valence in the task: participant learnt three  
325 labels (good-person, neutral-person, and bad-person) and three shapes (circle, square, and  
326 triangle), and the words related to identity, “self” or “other,” were presented in the shapes.  
327 As in 4a, participants were told to ignore the words inside the shape during the task.

328       Finally, experiment 5 was design to test the specificity of the moral valence. We  
329 extended experiment 1a with an additional independent variable: domains of the valence  
330 words. More specifically, besides the moral valence, we also added valence from other  
331 domains: appearance of person (beautiful, neutral, ugly), appearance of a scene (beautiful,

332 neutral, ugly), and emotion (happy, neutral, and sad). Label-shape pairs from different  
333 domains were separated into different blocks.

334 E-prime 2.0 was used for presenting stimuli and collecting behavioral responses,  
335 except that experiment 7a and 7b used Matlab Psychtoolbox (Brainard, 1997; Pelli, 1997).  
336 For participants recruited in Tsinghua University, they finished the experiment individually  
337 in a dim-lighted chamber, stimuli were presented on 22-inch CRT monitors and their head  
338 were fixed by a chin-rest brace. The distance between participants' eyes and the screen was  
339 about 60 cm. The visual angle of geometric shapes was about  $3.7^\circ \times 3.7^\circ$ , the fixation cross  
340 is of ( $0.8^\circ \times 0.8^\circ$  of visual angle) at the center of the screen. The words were of  $3.6^\circ \times 1.6^\circ$   
341 visual angle. The distance between the center of the shape or the word and the fixation  
342 cross was  $3.5^\circ$  of visual angle. For participants recruited in Wenzhou University, they  
343 finished the experiment in a group consisted of 3 ~ 12 participants in a dim-lighted testing  
344 room. Participants were required to finished the whole experiment independently. Also,  
345 they were instructed to start the experiment at the same time, so that the distraction  
346 between participants were minimized. The stimuli were presented on 19-inch CRT monitor.  
347 The visual angles are could not be exactly controlled because participants's chin were not  
348 fixed.

349 In most of these experiments, participant were also asked to fill a battery of  
350 questionnaire after they finish the behavioral tasks. All the questionnaire data are open  
351 (see, dataset 4 in Liu et al., 2020). See Table S1 for a summary information about all the  
352 experiments.

### 353 Data analysis

354 **Analysis of individual study.** We used the `tidyverse` of r (see script  
355 `Load_save_data.r`) to exclude the practicing trials, invalid trials of each participants, and  
356 invalid participants, if there were any, in the raw data. Results of each experiment were

357 then analyzed in two Bayesian approaches.

358        ***Bayesian hierarchical generalized linear model (BGLM).***

359        We first tested the effect of experimental manipulation using Bayesian hierarchical  
 360 generalized linear model (BGLM), because it provided three advantages over the classic  
 361 NHST approach (repeated measure ANOVA or t-tests): first, Bayesian models use  
 362 posterior distribution of parameter for statistical inference, therefore provided uncertainty  
 363 in estimation (Rouder & Lu, 2005), second, BGLM can use distribution that fit the real  
 364 distribution, which is the case for reaction time data (Rousselet & Wilcox, 2019), third,  
 365 BGLM also integrated different levels of analysis, fully account the variability from each  
 366 participants. We used the r package **BRMs** (Bürkner, 2017) to build the model, which used  
 367 Stan (Carpenter et al., 2017) to sample from the posterior.

368        ***Signal detection theory.***

369        As in (Hu, Lan, Macrae, & Sui, 2020; Sui, He, & Humphreys, 2012), we also used  
 370 signal detection approach to analyze the accuracy data. More specifically, we assume the  
 371 match trials are signal and the non-match trials are noise. To estimate the sensitivity and  
 372 criterion of SDT, we adopted the Bayesian hierarchical GLM approach from (Rouder & Lu,  
 373 2005). When modelling the accuracy data for one participant, we assume that the accuracy  
 374 of each trial is Bernoulli distributed (binomial with 1 trial), with probability  $p_i$  that  $y_i = 1$ .

$$y_i \sim \text{Bernoulli}(p_i)$$

375 In the perceptual matching task, the probability  $p_i$  can then be modeled as a function of  
 376 the trial type:

$$\Phi(p_i) = \beta_0 + \beta_1 \text{IsMatch}_i * \text{Valence}_i$$

377 The outcomes  $y_i$  are 0 if the participant responded “nonmatch” on trial  $i$ , 1 if they  
 378 responded “match.” The probability of the “match” response for trial  $i$  for a participant is

<sup>379</sup>  $p_i$ . We then write the generalized linear model on the probits (z-scores;  $\Phi$ , “Phi”) of  $ps$ .  $\Phi$   
<sup>380</sup> is the cumulative normal density function and maps  $z$  scores to probabilities. Given this  
<sup>381</sup> parameterization, the intercept of the model ( $\beta_0$ ) is the standardized false alarm rate  
<sup>382</sup> (probability of saying 1 when predictor is 0), which we take as our criterion  $c$ . The slope of  
<sup>383</sup> the model ( $\beta_1$ ) is the increase of saying 1 when predictor is 1, in  $z$ -scores, which is another  
<sup>384</sup> expression of  $d'$ . Therefore,  $c = -zHR = -\beta_0$ , and  $d' = \beta_1$ .

<sup>385</sup> In each experiment, we had multiple participants, to estimate the group-level  
<sup>386</sup> parameters, we need to estimate parameters on individual level and the group level  
<sup>387</sup> parameter simultaneously. In this case, as above, we first assume that the outcome of each  
<sup>388</sup> trial is Bernoulli distributed, with probability  $p_{ij}$  that  $y_{ij} = 1$ .

$$y_{ij} \sim Bernoulli(p_{ij})$$

<sup>389</sup> And the the generalized linear model was re-written to include two levels:

$$\Phi(p_{ij}) = \beta_{0j} + \beta_{1j} IsMatch_{ij} * Valence_{ij}$$

<sup>390</sup> The outcomes  $y_{ij}$  are 0 if participant  $j$  responded “nonmatch” on trial  $i$ , 1 if they  
<sup>391</sup> responded “match.” The probability of the “match” response for trial  $i$  for subject  $j$  is  $p_{ij}$ .  
<sup>392</sup> We again can write the generalized linear model on the probits (z-scores;  $\Phi$ , “Phi”) of  $ps$ .

<sup>393</sup> The subjective-specific intercepts ( $\beta_0 = -zFAR$ ) and slopes ( $\beta_1 = d'$ ) are describe  
<sup>394</sup> by multivariate normal with means and a covariance matrix for the parameters.

$$\begin{bmatrix} \beta_{0j} \\ \beta_{1j} \end{bmatrix} \sim N\left(\begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix}, \Sigma\right)$$

<sup>395</sup> For experiments that had 2 (matching: match vs. non-match) by 3 (moral character:  
<sup>396</sup> good vs. neutral vs. bad), i.e., experiment 1a, 1b, 1c, 2, 5, and 6a, the formula for BGLM is  
<sup>397</sup> as follow:

```

398     saymatch ~ 0 + Valence + Valence:ismatch + (0 + Valence +
399     Valence:ismatch | Subject), family = bernoulli(link="probit")

```

400       For experiments that had two by two by three design, we used the follow formula for  
 401 the BGLM:

```

402     saymatch ~ 0 + ID:Valence + ID:Valence:ismatch + (0 + ID:Valence +
403     ID:Valence:ismatch | Subject), family = bernoulli(link="probit")

```

404       For the reaction time, we used the log normal distribution  
 405 ([https://lindeloev.github.io/shiny-rt/#34\\_\(shifted\)\\_log-normal](https://lindeloev.github.io/shiny-rt/#34_(shifted)_log-normal)) to model the data. This  
 406 means that we need to estimate the posterior of two parameters:  $\mu$ ,  $\sigma$ .  $\mu$  is the mean of the  
 407 logNormal distribution, and  $\sigma$  is the disperse of the distribution. The log normal  
 408 distribution can be extended to shifted log normal distribution, with one more parameter:  
 409 shift, which is the earliest possible response. The reaction time is a linear function of trial  
 410 type:

$$y_{ij} = \beta_{0j} + \beta_{1j} * IsMatch_{ij} * Valence_{ij}$$

411       while the log of the reaction time is log-normal distributed:

$$\log(y_{ij}) \sim N(\mu_j, \sigma_j)$$

412  $y_{ij}$  is the RT of the  $i$ th trial of the  $j$ th participants.

$$\mu_j \sim N(\mu, \sigma)$$

$$\sigma_j \sim Cauchy()$$

413 Formula used for modeling the data as follow:

```

414     RT_sec ~ Valence*ismatch + (Valence*ismatch | Subject), family =
415     shifted_lognormal()

416     or

417     RT_sec ~ ID*Valence*ismatch + (ID*Valence*ismatch | Subject), family =
418     shifted_lognormal()

```

419 ***Hierarchical drift diffusion model (HDDM).***

420 To further explore the psychological mechanism under perceptual decision-making,  
 421 we used a generative mode drift diffusion model (DDM) to model our RTs and accuracy  
 422 data. As the hypothesis testing part, we also used hierarchical Bayesian model to fit the  
 423 DDM. The package we used was the HDDM (Wiecki, Sofer, & Frank, 2013), a python  
 424 package for fitting hierarchical DDM. We used the prior implemented in HDDM, that is,  
 425 weakly informative priors that constrains parameter estimates to be in the range of  
 426 plausible values based on past literature (Matzke & Wagenmakers, 2009). As reported in  
 427 Hu, Lan, Macrae, and Sui (2020), we used the stimulus code approach, match response  
 428 were coded as 1 and nonmatch responses were coded as 0. To fully explore all parameters,  
 429 we allow all four parameters of DDM free to vary. We then extracted the estimation of all  
 430 the four parameters for each participants for the correlation analyses. However, because  
 431 the starting point is only related to response (match vs. non-match) but not the valence of  
 432 the stimuli, we didn't included it in correlation analysis.

433 **Synthesized results.** Given that multiple experiments in the current study shared  
 434 similar experimental designs, We also synthesized their results to get a more precise and  
 435 robust estimation of the effect.

436 We used Bayesian hierarchical GLM model to synthesize the effect across different  
 437 studies by extending two-level hierarchical model into a three-level model, which  
 438 experiment as an additional level. For SDT, we can use a nested hierarchical model to

<sup>439</sup> model all the experiment with similar design:

$$y_{ijk} \sim Bernoulli(p_{ijk})$$

<sup>440</sup> where

$$\Phi(p_{ijk}) = \beta_{0jk} + \beta_{1jk} IsMatch_{ijk}$$

<sup>441</sup> The outcomes  $y_{ijk}$  are 0 if participant  $j$  in experiment  $k$  responded “nonmatch” on trial  $i$ ,

<sup>442</sup> 1 if they responded “match.”

$$\begin{bmatrix} \beta_{0jk} \\ \beta_{1jk} \end{bmatrix} \sim N\left(\begin{bmatrix} \theta_{0k} \\ \theta_{1k} \end{bmatrix}, \sum\right)$$

<sup>443</sup> and the experiment level parameter  $\mu_{0k}$  and  $\mu_{1k}$  is from a higher order

<sup>444</sup> distribution:

$$\begin{bmatrix} \theta_{0k} \\ \theta_{1k} \end{bmatrix} \sim N\left(\begin{bmatrix} \mu_0 \\ \mu_1 \end{bmatrix}, \sum\right)$$

<sup>445</sup> in which  $\mu_0$  and  $\mu_1$  means the population level parameter.

<sup>446</sup> In similar way, we expanded the RT model three-level model in which participants

<sup>447</sup> and experiments are two group level variable and participants were nested in the

<sup>448</sup> experiments.

$$\log(y_{ijk}) \sim N(\mu_{jk}, \sigma_{jk})$$

<sup>449</sup>  $y_{ijk}$  is the RT of the  $i$ th trial of the  $j$ th participants in the  $k$ th experiment.

$$\mu_{jk} \sim N(\mu_k, \sigma_k)$$

<sup>450</sup>

$$\sigma_{jk} \sim Cauchy()$$

451

$$\mu_k \sim N(\mu, \sigma)$$

452

$$\theta_k \sim Cauchy()$$

453        Using the Bayesian hierarchical model, we can directly estimate the over-all effect of  
 454        valence on  $d'$  and RT across all experiments with similar experimental design, instead of  
 455        using a two-step approach where we first estimate the  $d'$  for each participant and then use  
 456        a random effect model meta-analysis (Goh, Hall, & Rosenthal, 2016).

457        *Effect of moral character.*

458        We synthesized effect size of  $d'$  and RT from experiment 1a, 1b, 1c, 2, 5, and 6a for  
 459        the effect of moral character. We reported the synthesized the effect across all experiments  
 460        that tested the valence effect, using the mini meta-analysis approach (Goh, Hall, &  
 461        Rosenthal, 2016).

462        ***Effect of moral self.***

463        We further synthesized the effect of moral self, which included results from  
 464        experiment 3a, 3b, and 6b. In these experiment, we directly tested two possible  
 465        explanations: moral self as social categorization process and value-based attention.

466        ***Implicit interaction between valence and self-relevance.***

467        In the third part, we focused on experiment 4a and 4b, which were designed to  
 468        examine two more nuanced explanation concerning the good-self. The design of experiment  
 469        4a and 4b are complementary. Together, they can test whether participants are more  
 470        sensitive to the moral character of the Self (4a), or the identity of the morally Good (4b).

471        ***Specificity of the valence effect.***

472        In this part, we reported the data from experiment 5, which included positive,  
 473        neutral, and negative valence from four different domains: morality, aesthetic of person,

474 aesthetic of scene, and emotion. This experiment was design to test whether the positive  
475 bias is specific to morality.

476 ***Behavior-Questionnaire correlation.***

477 Finally, we explored correlation between results from behavioral results and  
478 self-reported measures.

479 For the questionnaire part, we are most interested in the self-rated distance between  
480 different person and self-evaluation related questionnaires: self-esteem, moral-self identity,  
481 and moral self-image. Other questionnaires (e.g., personality) were not planned to  
482 correlated with behavioral data were not included. Note that all questionnaire data were  
483 reported in (Liu et al., 2020).

484 For the behavioral task part, we used three parameters from drift diffusion model:  
485 drift rate ( $v$ ), boundary separation ( $a$ ), and non decision-making time ( $t$ ), because these  
486 parameters has relative clear psychological meaning. We used the mean of parameter  
487 posterior distribution as the estimate of each parameter for each participants in the  
488 correlation analysis. We used alpha = 0.05 and used bootstrap by BootES package (Kirby  
489 & Gerlanc, 2013) to estimate the correlation.

490 **Part 1: Perceptual processing moral character related inforation**

491 In this part, we report five experiments that tested whether an associative learning  
492 task, in which concepts of moral character are associated with geometric shapes, will  
493 impact the perceptual decision-making.

494 **Experiment 1a**

495 **Methods.**

**496      *Participants.***

497      57 college students (38 female, age =  $20.75 \pm 2.54$  years) participated. 39 of them  
498      were recruited from Tsinghua University community in 2014; 18 were recruited from  
499      Wenzhou University in 2017. All participants were right-handed except one, and all had  
500      normal or corrected-to-normal vision. Informed consent was obtained from all participants  
501      prior to the experiment according to procedures approved by the local ethics committees. 6  
502      participant's data were excluded from analysis because nearly random level of accuracy,  
503      leaving 51 participants (34 female, age =  $20.72 \pm 2.44$  years).

**504      *Stimuli and Tasks.***

505      Three geometric shapes were used in this experiment: triangle, square, and circle.  
506      These shapes were paired with three labels (bad person, good person or neutral person).  
507      The pairs were counterbalanced across participants.

**508      *Procedure.***

509      This experiment had two phases. First, there was a brief learning stage. Participants  
510      were asked to learn the relationship between geometric shapes (triangle, square, and circle)  
511      and different concepts of moral character (bad person, a good person, or a neutral person).  
512      For example, a participant was told, "bad person is a circle; good person is a triangle; and  
513      a neutral person is a square." After participants remembered the associations (usually in a  
514      few minutes), they started a practicing phase of matching task which had the exact task as  
515      in the experimental task.

516      In the experimental task, participants judged whether shape-label pairs, which were  
517      subsequently presented, were correct (i.e., the same as they learned). Each trial started  
518      with the presentation of a central fixation cross for 500 ms. Subsequently, a pairing of a  
519      shape and label (good person, bad person, and neutral person) was presented for 100 ms.  
520      The pair presented could confirm to the verbal instruction for each pairing given in the  
521      training stage, or it could be a recombination of a shape with a different label, with the

shape-label pairings being generated at random. The next frame showed a blank for 1100ms. Participants were expected to judge whether the shape was correctly assigned to the person by pressing one of the two response buttons as quickly and accurately as possible within this timeframe (to encourage immediate responding). Feedback (correct or incorrect) was given on the screen for 500 ms at the end of each trial, if no response detected, “too slow” was presented to remind participants to accelerate. Participants were informed of their overall accuracy at the end of each block. The practice phase finished and the experimental task began after the overall performance of accuracy during practice phase achieved 60%.

For participants from the Tsinghua community, they completed 6 experimental blocks of 60 trials. Thus, there were 60 trials in each condition (bad-person match, bad-person nonmatch, good-person match, good-person nonmatch, neutral-person match, and neutral-person nonmatch). For the participants from Wenzhou University, they finished 6 blocks of 120 trials, therefore, 120 trials for each condition.

### 536 ***Data analysis.***

As described in general methods section, we used Bayesian Bayesian Hierarchical Generalized Linear Model for hypothesis testing and Hierarchical drift diffusion model. We also included the classic NHST results in the online supplementary results.

### 540 **Results.**

#### 541 ***Hypothesis testing.***

542 *d prime.*

We fitted a Bayesian hierarchical GLM for signal detection theory. The results showed that when the shapes were tagged with labels with different moral character, the sensitivity ( $d'$ ) and criteria ( $c$ ) were both influenced. For the  $d'$ , we found that the shapes associated with good person (2.46, 95% CI[2.21 2.72]) is greater than shapes tagged with moral bad (2.07, 95% CI[1.83 2.32]),  $P_{PosteriorComparison} = 1$ . Shape tagged with morally

548 good person is also greater than shapes tagged with neutral person (2.23, 95% CI[1.95  
 549 2.49]),  $P_{PosteriorComparison} = 0.97$ . Also, the shapes tagged with neutral person is greater  
 550 than shapes tagged with morally bad person,  $P_{PosteriorComparison} = 0.92$ .

551 Interesting, we also found the criteria for three conditions also differ, the shapes  
 552 tagged with good person has the highest criteria (-1.01, [-1.14 -0.88]), followed by shapes  
 553 tagged with neutral person(-1.06, [-1.21 -0.92]), and then the shapes tagged with bad  
 554 person(-1.11, [-1.25 -0.97]). However, pair-wise comparison showed that only showed strong  
 555 evidence for the difference between good and bad conditions.

556 *Reaction times.*

557 We fitted a Bayesian hierarchical GLM for RTs, with a log-normal distribution as the  
 558 link function. We used the posterior distribution of the regression coefficient to make  
 559 statistical inferences. As in previous studies, the matched conditions are much faster than  
 560 the mismatched trials ( $P_{PosteriorComparison} = 1$ ). We focused on matched trials only, and  
 561 compared different conditions: Good is faster than the neutral,  $P_{PosteriorComparison} = .99$ ,  
 562 it was also faster than the Bad condition,  $P_{PosteriorComparison} = 1$ . And the neutral  
 563 condition is faster than the bad condition,  $P_{PosteriorComparison} = .99$ . However, the  
 564 mismatched trials are largely overlapped. See Figure ??.

565 **HDDM.**

566 We fitted our data with HDDM, using the response-coding (See also, Hu, Lan,  
 567 Macrae, & Sui, 2020). We estimated separate drift rate ( $v$ ), non-decision time ( $T_0$ ), and  
 568 boundary separation ( $a$ ) for each condition. We found that the shapes tagged with good  
 569 person has higher drift rate and higher boundary separation than shapes tagged with both  
 570 neutral and bad person. Also, the shapes tagged with neutral person has a higher drift rate  
 571 than shapes tagged with bad person, but not for the boundary separation. Finally, we  
 572 found that shapes tagged with bad person had longer non-decision time (see Figure ??).

573 **Experiment 1b**

574 This study was conducted to further confirm that the moral character information  
575 influence the perceptual decision making instead of other factors such as the familiarity of  
576 words. To do so, we selected different words whose dominant meaning is related to moral  
577 character but with similar level of familiarity between different words.

578 **Method.**

579 ***Participants.***

580 72 college students (49 female, age =  $20.17 \pm 2.08$  years) participated. 39 of them  
581 were recruited from Tsinghua University community in 2014; 33 were recruited from  
582 Wenzhou University in 2017. All participants were right-handed except one, and all had  
583 normal or corrected-to-normal vision. Informed consent was obtained from all participants  
584 prior to the experiment according to procedures approved by the local ethics committees.  
585 20 participant's data were excluded from analysis because nearly random level of accuracy,  
586 leaving 52 participants (36 female, age =  $20.25 \pm 2.31$  years).

587 **Stimuli and Tasks.** Three geometric shapes (triangle, square, and circle, with  $3.7^\circ$   
588  $\times 3.7^\circ$  of visual angle) were presented above a white fixation cross subtending  $0.8^\circ \times 0.8^\circ$   
589 of visual angle at the center of the screen. The three shapes were randomly assigned to  
590 three labels with different moral valence: a morally bad person (" , " ERen), a morally  
591 good person (" , " ShanRen) or a morally neutral person (" , " ChangRen). The order of  
592 the associations between shapes and labels was counterbalanced across participants.

593 Three labels used in this experiment was selected based on the rating results from an  
594 independent survey, in which participants rated the familiarity, frequency, and concreteness  
595 of eight different words online. Of the eight words, three of them are morally positive  
596 (HaoRen, ShanRen, Junzi), two of them are morally neutral (ChangRen, FanRen), and  
597 three of them are morally negative (HuaiRen, ERen, LiuMang). An independent sample  
598 consist of 35 participants (22 females, age  $20.6 \pm 3.11$ ) were recruited to rate these words.

599 Based on the ratings (see supplementary materials Figure S1), we selected ShanRen,  
600 ChangRen, and ERen to represent morally positive, neutral, and negative person.

601 ***Procedure.***

602 For participants from both Tsinghua community and Wenzhou community, the  
603 procedure in the current study was exactly same as in experiment 1a.

604 **Data Analysis.** Data was analyzed as in experiment 1a.

605 **Results.**

606 **NHST.**

607 Figure ?? shows  $d$  prime and reaction times of experiment 1b.

608  $d$  prime.

609 Repeated measures ANOVA revealed main effect of valence,  $F(1.83, 93.20) = 14.98$ ,  
610  $MSE = 0.18$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .053$ . Paired t test showed that the Good-Person condition  
611 ( $1.87 \pm 0.102$ ) was with greater  $d$  prime than Neutral condition ( $1.44 \pm 0.101$ ,  $t(51) =$   
612  $5.945$ ,  $p < 0.001$ ). We also found that the Bad-Person condition ( $1.67 \pm 0.11$ ) has also  
613 greater  $d$  prime than neutral condition ,  $t(51) = 3.132$ ,  $p = 0.008$ ). There Good-person  
614 condition was also slightly greater than the bad condition,  $t(51) = 2.265$ ,  $p = 0.0701$ .

615 *Reaction times.*

616 We found interaction between Matchness and Valence ( $F(1.95, 99.31) = 19.71$ ,  
617  $MSE = 960.92$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .031$ ) and then analyzed the matched trials and  
618 mismatched trials separately, as in experiment 1a. For matched trials, we found the effect  
619 of valence  $F(1.94, 99.10) = 33.97$ ,  $MSE = 1,343.19$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .115$ . Post-hoc  $t$ -tests  
620 revealed that shapes associated with Good Person ( $684 \pm 8.77$ ) were responded faster than  
621 Neutral-Person ( $740 \pm 9.84$ ), ( $t(51) = -8.167$ ,  $p < 0.001$ ) and Bad Person ( $728 \pm 9.15$ ),  
622  $t(51) = -5.724$ ,  $p < 0.0001$ ). While there was no significant differences between Neutral and

623 Bad-Person condition ( $t(51) = 1.686, p = 0.221$ ). For non-matched trials, there was no  
 624 significant effect of Valence ( $F(1.90, 97.13) = 1.80, MSE = 430.15, p = .173, \hat{\eta}_G^2 = .003$ ).

625 **BGLM.**

626 *Signal detection theory analysis of accuracy.*

627 We fitted a Bayesian hierarchical GLM for SDT. The results showed that when the  
 628 shapes were tagged with labels with different moral valence, the sensitivity ( $d'$ ) and criteria  
 629 ( $c$ ) were both influence. For the  $d'$ , we found that the shapes tagged with morally good  
 630 person (2.46, 95% CI[2.21 2.72]) is greater than shapes tagged with moral bad (2.07, 95%  
 631 CI[1.83 2.32]),  $P_{PosteriorComparison} = 1$ . Shape tagged with morally good person is also  
 632 greater than shapes tagged with neutral person (2.23, 95% CI[1.95 2.49]),  
 633  $P_{PosteriorComparison} = 0.97$ . Also, the shapes tagged with neutral person is greater than  
 634 shapes tagged with morally bad person,  $P_{PosteriorComparison} = 0.92$ .

635 Interesting, we also found the criteria for three conditions also differ, the shapes  
 636 tagged with good person has the highest criteria (-1.01, [-1.14 -0.88]), followed by shapes  
 637 tagged with neutral person(1.06, [-1.21 -0.92]), and then the shapes tagged with bad  
 638 person(-1.11, [-1.25 -0.97]). However, pair-wise comparison showed that only showed strong  
 639 evidence for the difference between good and bad conditions.

640 *Reaction time.*

641 We fitted a Bayesian hierarchical GLM for RTs, with a log-normal distribution as the  
 642 link function. We used the posterior distribution of the regression coefficient to make  
 643 statistical inferences. As in previous studies, the matched conditions are much faster than  
 644 the mismatched trials ( $P_{PosteriorComparison} = 1$ ). We focused on matched trials only, and  
 645 compared different conditions: Good is faster than the neutral,  $P_{PosteriorComparison} = .99$ ,  
 646 it was also faster than the Bad condition,  $P_{PosteriorComparison} = 1$ . And the neutral  
 647 condition is faster than the bad condition,  $P_{PosteriorComparison} = .99$ . However, the  
 648 mismatched trials are largely overlapped. See Figure ??.

**HDDM.**

We found that the shapes tagged with good person has higher drift rate and higher boundary separation than shapes tagged with both neutral and bad person. Also, the shapes tagged with neutral person has a higher drift rate than shapes tagged with bad person, but not for the boundary separation. Finally, we found that shapes tagged with bad person had longer non-decision time (see figure ??).

**Discussion.** These results confirmed the facilitation effect of positive moral valence on the perceptual matching task. This pattern of results mimic prior results demonstrating self-bias effect on perceptual matching (Sui, He, & Humphreys, 2012) and in line with previous studies that indirect learning of other's moral reputation do have influence on our subsequent behavior (Fouragnan et al., 2013).

**Experiment 1c**

In this study, we further control the valence of words using in our experiment.

Instead of using label with moral valence, we used valence-neutral names in China. Participant first learn behaviors of the different person, then, they associate the names and shapes. And then they perform a name-shape matching task.

**Method.*****Participants.***

23 college students (15 female, age =  $22.61 \pm 2.62$  years) participated. All of them were recruited from Tsinghua University community in 2014. Informed consent was obtained from all participants prior to the experiment according to procedures approved by the local ethics committees. No participant was excluded because they overall accuracy were above 0.6.

***Stimuli and Tasks.***

673 Three geometric shapes (triangle, square, and circle, with  $3.7^\circ \times 3.7^\circ$  of visual angle)

674 were presented above a white fixation cross subtending  $0.8^\circ \times 0.8^\circ$  of visual angle at the

675 center of the screen. The three most common names were chosen, which are neutral in

676 moral valence before the manipulation. Three names (Zhang, Wang, Li) were first paired

677 with three paragraphs of behavioral description. Each description includes one sentence of

678 biographic information and four sentences that describing the moral behavioral under that

679 name. To assess the that these three descriptions represented good, neutral, and bad

680 valence, we collected the ratings of three person on six dimensions: morality, likability,

681 trustworthiness, dominance, competence, and aggressiveness, from an independent sample

682 ( $n = 34$ , 18 female, age =  $19.6 \pm 2.05$ ). The rating results showed that the person with

683 morally good behavioral description has higher score on morality ( $M = 3.59$ ,  $SD = 0.66$ )

684 than neutral ( $M = 0.88$ ,  $SD = 1.1$ ),  $t(33) = 12.94$ ,  $p < .001$ , and bad conditions ( $M = -3.4$ ,

685  $SD = 1.1$ ),  $t(33) = 30.78$ ,  $p < .001$ . Neutral condition was also significant higher than bad

686 conditions  $t(33) = 13.9$ ,  $p < .001$  (See supplementary materials).

687 ***Procedure.***

688 After arriving the lab, participants were informed to complete two experimental

689 tasks, first a social memory task to remember three person and their behaviors, after tested

690 for their memory, they will finish a perceptual matching task. In the social memory task,

691 the descriptions of three person were presented without time limitation. Participant

692 self-paced to memorized the behaviors of each person. After they memorizing, a

693 recognition task was used to test their memory effect. Each participant was required to

694 have over 95% accuracy before preceding to matching task. The perceptual learning task

695 was followed, three names were randomly paired with geometric shapes. Participants were

696 required to learn the association and perform a practicing task before they start the formal

697 experimental blocks. They kept practicing until they reached 70% accuracy. Then, they

698 would start the perceptual matching task as in experiment 1a. They finished 6 blocks of

699 perceptual matching trials, each have 120 trials.

700       **Data Analysis.** Data was analyzed as in experiment 1a.

701       **Results.** Figure ?? shows  $d$  prime and reaction times of experiment 1c. We  
 702       conducted same analysis as in Experiment 1a. Our analysis didn't show effect of valence  
 703       on  $d$  prime,  $F(1.93, 42.56) = 0.23$ ,  $MSE = 0.41$ ,  $p = .791$ ,  $\hat{\eta}_G^2 = .005$ . Neither the effect of  
 704       valence on RT ( $F(1.63, 35.81) = 0.22$ ,  $MSE = 2,212.71$ ,  $p = .761$ ,  $\hat{\eta}_G^2 = .001$ ) or  
 705       interaction between valence and matchness on RT ( $F(1.79, 39.43) = 1.20$ ,  
 706        $MSE = 1,973.91$ ,  $p = .308$ ,  $\hat{\eta}_G^2 = .005$ ).

707       ***Signal detection theory analysis of accuracy.***

708       We fitted a Bayesian hierarchical GLM for SDT. The results showed that when the  
 709       shapes were tagged with labels with different moral valence, the sensitivity ( $d'$ ) and criteria  
 710       ( $c$ ) were both influenced. For the  $d'$ , we found that the shapes tagged with morally good  
 711       person (2.30, 95% CI[1.93 2.70]) is greater than shapes tagged with moral bad (2.11, 95%  
 712       CI[1.83 2.42]),  $P_{PosteriorComparison} = 0.8$ . Shape tagged with morally good person is also  
 713       greater than shapes tagged with neutral person (2.16, 95% CI[1.88 2.45]),  
 714        $P_{PosteriorComparison} = 0.75$ .

715       Interestingly, we also found the criteria for three conditions also differ, the shapes  
 716       tagged with good person has the highest criteria (-0.97, [-1.12 -0.82]), followed by shapes  
 717       tagged with neutral person(-0.96, [-1.09 -0.83]), and then the shapes tagged with bad  
 718       person(-1.03, [-1.22 -0.84]). However, pair-wise comparison showed that only showed strong  
 719       evidence for the difference between good and bad conditions.

720       ***Reaction time.***

721       We fitted a Bayesian hierarchical GLM for RTs, with a log-normal distribution as the  
 722       link function. We used the posterior distribution of the regression coefficient to make  
 723       statistical inferences. As in previous studies, the matched conditions are much faster than  
 724       the mismatched trials ( $P_{PosteriorComparison} = .75$ ). We focused on matched trials only, and  
 725       compared different conditions: Good () is not faster than the neutral (),

<sup>726</sup>  $P_{PosteriorComparison} = .5$ , it was faster than the Bad condition (),

<sup>727</sup>  $P_{PosteriorComparison} = .88$ . And the neutral condition is faster than the bad condition,

<sup>728</sup>  $P_{PosteriorComparison} = .95$ . However, the mismatched trials are largely overlapped.

<sup>729</sup> **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu,  
<sup>730</sup> Lan, Macrae, & Sui, 2020). We estimated separate drift rate ( $v$ ), non-decision time ( $T_0$ ),  
<sup>731</sup> and boundary separation ( $a$ ) for each condition. We found that the shapes tagged with  
<sup>732</sup> good person has higher drift rate and higher boundary separation than shapes tagged with  
<sup>733</sup> both neutral and bad person. Also, the shapes tagged with neutral person has a higher drift  
<sup>734</sup> rate than shapes tagged with bad person, but not for the boundary separation. Finally, we  
<sup>735</sup> found that shapes tagged with bad person had longer non-decision time (see figure ??)).

<sup>736</sup> **Experiment 2: Sequential presenting**

<sup>737</sup> Experiment 2 was conducted for two purpose: (1) to further confirm the facilitation  
<sup>738</sup> effect of positive moral associations; (2) to test the effect of expectation of occurrence of  
<sup>739</sup> each pair. In this experiment, after participant learned the association between labels and  
<sup>740</sup> shapes, they were presented a label first and then a shape, they then asked to judge  
<sup>741</sup> whether the shape matched the label or not (see (Sui, Sun, Peng, & Humphreys, 2014)).  
<sup>742</sup> Previous studies showed that when the labels presented before the shapes, participants  
<sup>743</sup> formed expectations about the shape, and therefore a top-down process were introduced  
<sup>744</sup> into the perceptual matching processing. If the facilitation effect of positive moral valence  
<sup>745</sup> we found in experiment 1 was mainly drive by top-down processes, this sequential  
<sup>746</sup> presenting paradigm may eliminate or attenuate this effect; if, however, the facilitation  
<sup>747</sup> effect occurred because of button-up processes, then, similar facilitation effect will appear  
<sup>748</sup> even with sequential presenting paradigm.

<sup>749</sup> **Method.**

<sup>750</sup> **Participants.**

751        35 participants (17 female, age =  $21.66 \pm 3.03$ ) were recruited. 24 of them had  
752        participated in Experiment 1a (9 male, mean age = 21.9, s.d. = 2.9), and the time gap  
753        between these experiment 1a and experiment 2 is at least six weeks. The results of 1  
754        participants were excluded from analysis because of less than 60% overall accuracy,  
755        remains 34 participants (17 female, age =  $21.74 \pm 3.04$ ).

756        ***Procedure.***

757        In Experiment 2, the sequential presenting makes the matching task much easier than  
758        experiment 1. To avoid ceiling effect on behavioral data, we did a few pilot experiments to  
759        get optimal parameters, i.e., the conditions under which participant have similar accuracy  
760        as in Experiment 1 (around 70 ~ 80% accuracy). In the final procedure, the label (good  
761        person, bad person, or neutral person) was presented for 50 ms and then masked by a  
762        scrambled image for 200 ms. A geometric shape followed the scrambled mask for 50 ms in  
763        a noisy background (which was produced by first decomposing a square with  $\frac{3}{4}$  gray area  
764        and  $\frac{1}{4}$  white area to small squares with a size of  $2 \times 2$  pixels and then re-combine these  
765        small pieces randomly), instead of pure gray background in Experiment 1. After that, a  
766        blank screen was presented 1100 ms, during which participants should press a button to  
767        indicate the label and the shape match the original association or not. Feedback was given,  
768        as in study 1. The next trial then started after 700 ~ 1100 ms blank. Other aspects of  
769        study 2 were identical to study 1.

770        ***Data analysis.***

771        Data was analyzed as in study 1a.

772        **Results.**

773        **NHST.**

774        Figure ?? shows  $d'$  prime and reaction times of experiment 2. Less than 0.2% correct  
775        trials with less than 200ms reaction times were excluded.

776 *d prime.*

777 There was evidence for the main effect of valence,  $F(1.83, 60.36) = 14.41$ ,

778  $MSE = 0.23$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .066$ . Paired t test showed that the Good-Person condition

779 ( $2.79 \pm 0.17$ ) was with greater *d* prime than Netural condition ( $2.21 \pm 0.16$ ,  $t(33) = 4.723$ ,

780  $p = 0.001$ ) and Bad-person condition ( $2.41 \pm 0.14$ ),  $t(33) = 4.067$ ,  $p = 0.008$ ). There was

781 no-significant difference between Neutral-person and Bad-person condition,  $t(33) = -1.802$ ,

782  $p = 0.185$ .

783 *Reaction time.*

784 The results of reaction times of matchness trials showed similar pattern as the *d*

785 prime data.

786 We found interaction between Matchness and Valence ( $F(1.99, 65.70) = 9.53$ ,

787  $MSE = 605.36$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .017$ ) and then analyzed the matched trials and

788 mismatched trials separately, as in experiment 1a. For matched trials, we found the effect

789 of valence  $F(1.99, 65.76) = 10.57$ ,  $MSE = 1,192.65$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .067$ . Post-hoc *t*-tests

790 revealed that shapes associated with Good Person ( $548 \pm 9.4$ ) were responded faster than

791 Neutral-Person ( $582 \pm 10.9$ ), ( $t(33) = -3.95$ ,  $p = 0.0011$ ) and Bad Person ( $582 \pm 10.2$ ),

792  $t(33) = -3.9$ ,  $p = 0.0013$ ). While there was no significant differences between Neutral and

793 Bad-Person condition ( $t(33) = -0.01$ ,  $p = 0.999$ ). For non-matched trials, there was no

794 significant effect of Valence ( $F(1.99, 65.83) = 0.17$ ,  $MSE = 489.80$ ,  $p = .843$ ,  $\hat{\eta}_G^2 = .001$ ).

795 **BGLMM.**

796 *Signal detection theory analysis of accuracy.*

797 We fitted a Bayesian hierarchical GLM for SDT. The results showed that when the

798 shapes were tagged with labels with different moral valence, the sensitivity ( $d'$ ) and criteria

799 ( $c$ ) were both influence. For the  $d'$ , we found that the shapes tagged with morally good

800 person ( $2.46$ , 95% CI[ $2.21$   $2.72$ ]) is greater than shapes tagged with moral bad ( $2.07$ , 95%

801 CI[ $1.83$   $2.32$ ]),  $P_{PosteriorComparison} = 1$ . Shape tagged with morally good person is also

802 greater than shapes tagged with neutral person (2.23, 95% CI[1.95 2.49]),  
 803  $P_{PosteriorComparison} = 0.97$ . Also, the shapes tagged with neutral person is greater than  
 804 shapes tagged with morally bad person,  $P_{PosteriorComparison} = 0.92$ .

805 Interesting, we also found the criteria for three conditions also differ, the shapes  
 806 tagged with good person has the highest criteria (-1.01, [-1.14 -0.88]), followed by shapes  
 807 tagged with neutral person(1.06, [-1.21 -0.92]), and then the shapes tagged with bad  
 808 person(-1.11, [-1.25 -0.97]). However, pair-wise comparison showed that only showed strong  
 809 evidence for the difference between good and bad conditions.

810 *Reaction times.*

811 We fitted a Bayesian hierarchical GLM for RTs, with a log-normal distribution as the  
 812 link function. We used the posterior distribution of the regression coefficient to make  
 813 statistical inferences. As in previous studies, the matched conditions are much faster than  
 814 the mismatched trials ( $P_{PosteriorComparison} = .75$ ). We focused on matched trials only, and  
 815 compared different conditions: Good () is not faster than the neutral (),  
 816  $P_{PosteriorComparison} = .5$ , it was faster than the Bad condition (),  
 817  $P_{PosteriorComparison} = .88$ . And the neutral condition is faster than the bad condition,  
 818  $P_{PosteriorComparison} = .95$ . However, the mismatched trials are largely overlapped.

819 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu,  
 820 Lan, Macrae, & Sui, 2020). We estimated separate drift rate ( $v$ ), non-decision time ( $T_0$ ),  
 821 and boundary separation ( $a$ ) for each condition. We found that the shapes tagged with  
 822 good person has higher drift rate and higher boundary separation than shapes tagged with  
 823 both neutral and bad person. Also, the shapes tagged with neutral person has a higher  
 824 drift rate than shapes tagged with bad person, but not for the boundary separation.  
 825 Finally, we found that shapes tagged with bad person had longer non-decision time (see  
 826 figure @ref(fig:plot-exp1c -HDDM)).

827 **Discussion**

828 In this experiment, we repeated the results pattern that the positive moral valenced  
829 stimuli has an advantage over the neutral or the negative valence association. Moreover,  
830 with a cross-task analysis, we did not find evidence that the experiment task interacted  
831 with moral valence, suggesting that the effect might not be effect by experiment task.  
832 These findings suggested that the facilitation effect of positive moral valence is robust and  
833 not affected by task. This robust effect detected by the associative learning is unexpected.

834 **Experiment 6a: EEG study 1**

835 Experiment 6a was conducted to study the neural correlates of the positive  
836 prioritization effect. The behavioral paradigm is same as experiment 2.

837 **Method.**

838 ***Participants.***

839 24 college students (8 female, age =  $22.88 \pm 2.79$ ) participated the current study, all  
840 of them were from Tsinghua University in 2014. Informed consent was obtained from all  
841 participants prior to the experiment according to procedures approved by a local ethics  
842 committee. No participant was excluded from behavioral analysis.

843 **Experimental design.** The experimental design of this experiment is same as  
844 experiment 2: a  $3 \times 2$  within-subject design with moral valence (good, neutral and bad  
845 associations) and matchness between shape and label (match vs. mismatch for the personal  
846 association) as within-subject variables.

847 ***Stimuli.***

848 Three geometric shapes (triangle, square and circle, each  $4.6^\circ \times 4.6^\circ$  of visual angle)  
849 were presented at the center of screen for 50 ms after 500ms of fixation ( $0.8^\circ \times 0.8^\circ$  of  
850 visual angle). The association of the three shapes to bad person (“ , HuaiRen”), good

851 person (“ , HaoRen”) or ordinary person (“ , ChangRen”) was counterbalanced across  
852 participants. The words bad person, good person or ordinary person ( $3.6^\circ \times 1.6^\circ$ ) was also  
853 displayed at the center fo the screen. Participants had to judge whether the pairings of  
854 label and shape matched (e.g., Does the circle represent a bad person?). The experiment  
855 was run on a PC using E-prime software (version 2.0). These stimuli were displayed on a  
856 22-in CRT monitor ( $1024 \times 768$  at 100Hz). We used backward masking to avoid  
857 over-processing of the moral words, in which a scrambled picture were presented for 900 ms  
858 after the label. Also, to avoid the ceiling effect on accuracy, shapes were presented on a  
859 noisy background based on our pilot studies. The noisy images were made by scrambling a  
860 picture of 3/4gray and 1/4 white at resolution of  $2 \times 2$  pixel.

861 ***Procedure.***

862 The procedure was similar to Experiment 2. Participants finished 9 blocks of trial,  
863 each with 120 trials. In total, participants finished 180 trials for each combination of  
864 condition.

865 As in experiment 2 (Sui, He, & Humphreys, 2012), subjects first learned the  
866 associations between labels and shapes and then completed a shape-label matching task  
867 (e.g., good person-triangle). In each trial of the matching task, a fixation were first  
868 presented for 500 ms, followed by a 50 ms label; then, a scrambled picture presented 900  
869 ms. After the backward mask, the shape were presented on a noisy background for 50ms.  
870 Participant have to response in 1000ms after the presentation of the shape, and finally, a  
871 feedback screen was presented for 500 ms (see figure 1). The inter-trial interval (ITI) were  
872 randomly varied at the range of 1000 ~ 1400 ms.

873 All the stimuli were presented on a gray background (RGB: 127, 127, 127). E-primed  
874 2.0 was used to present stimuli and collect behavioral results. Data were collected and  
875 analyzed when accuracy performance in total reached 60%.

876 **Data Analysis.** Data was analyzed as in experiment 1a.

**Results.****NHST.**

Only the behavioral results were reported here. Figure ?? shows  $d$  prime and reaction times of experiment 6a.

$d$  prime.

We conducted repeated measures ANOVA, with moral valence as independent variable. The results revealed the main effect of valence ( $F(1.74, 40.05) = 3.76$ ,  $MSE = 0.10$ ,  $p = .037$ ,  $\hat{\eta}_G^2 = .021$ ). Post-hoc analysis revealed that shapes link with Good person (mean = 3.13, SE = 0.109) is greater than Neutral condition (mean = 2.88, SE = 0.14),  $t = 2.916$ ,  $df = 24$ ,  $p = 0.02$ , p-value adjusted by Tukey method, but the  $d$  prime between Good and bad (mean = 3.03, SE = 0.142) ( $t = 1.512$ ,  $df = 24$ ,  $p = 0.3034$ , p-value adjusted by Tukey method), bad and neutral ( $t = 1.599$ ,  $df = 24$ ,  $p = 0.2655$ , p-value adjusted by Tukey method) were not significant.

*Reaction times.*

The results of reaction times of matchness trials showed similar pattern as the  $d$  prime data.

We found intercation between Matchness and Valence ( $F(1.97, 45.20) = 20.45$ ,  $MSE = 450.47$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .021$ ) and then analyzed the matched trials and mismatched trials separately, as in experiment 2. For matched trials, we found the effect of valence  $F(1.97, 45.25) = 32.37$ ,  $MSE = 522.42$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .078$ . For non-matched trials, there was no significant effect of Valence ( $F(1.77, 40.67) = 0.35$ ,  $MSE = 242.15$ ,  $p = .679$ ,  $\hat{\eta}_G^2 = .000$ ). Post-hoc  $t$ -tests revealed that shapes associated with Good Person (mean = 550, SE = 13.8) were responded faster than Neutral-Person (501, SE = 14.7), ( $t(24) = -5.171$ ,  $p = 0.0001$ ) and Bad Person (523, SE = 16.3),  $t(24) = -8.137$ ,  $p < 0.0001$ ., and Neutral is faster than Bad-Person condition ( $t(32) = -3.282$ ,  $p = 0.0085$ ).

**BGLM.**

903        *Signal detection theory analysis of accuracy.*

904        *Reaction time.*

905        **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu,  
906 Lan, Macrae, & Sui, 2020). We estimated separate drift rate ( $v$ ), non-decision time ( $T_0$ ),  
907 and boundary separation ( $a$ ) for each condition. We found that, similar to experiment 2,  
908 the shapes tagged with good person has higher drift rate and higher boundary separation  
909 than shapes tagged with both neutral and bad person, but only for the self-referential  
910 condition. Also, the shapes tagged with neutral person has a higher drift rate than shapes  
911 tagged with bad person, but not for the boundary separation, and this effect also exist only  
912 for the self-referential condition.

913        Interestingly, we found that in both self-referential and other-referential conditions,  
914 the shapes associated bad valence have higher drift rate and higher boundary separation.  
915 which might suggest that the shape associated with bad stimuli might be prioritized in the  
916 non-match trials (see figure ??).

917        **Part 2: interaction between valence and identity**

918        In this part, we report two experiments that aimed at testing whether the moral  
919 valence effect found in the previous experiment can be modulated by the self-referential  
920 processing.

921        RUE## Experiment 3a To examine the modulation effect of positive valence was an  
922 intrinsic, self-referential process, we designed study 3. In this study, moral valence was  
923 assigned to both self and a stranger. We hypothesized that the modulation effect of moral  
924 valence will be stronger for the self than for a stranger.

925        **Method.**

926        **Participants.**

927        38 college students (15 female, age =  $21.92 \pm 2.16$ ) participated in experiment 3a.

928        All of them were right-handed, and all had normal or corrected-to-normal vision. Informed  
929        consent was obtained from all participants prior to the experiment according to procedures  
930        approved by a local ethics committee. One female and one male student did not finish the  
931        experiment, and 1 participants' data were excluded from analysis because less than 60%  
932        overall accuracy, remains 35 participants (13 female, age =  $22.11 \pm 2.13$ ).

933        ***Design.***

934        Study 3a combined moral valence with self-relevance, hence the experiment has a  $2 \times$   
935         $3 \times 2$  within-subject design. The first variable was self-relevance, include two levels:  
936        self-relevance vs. stranger-relevance; the second variable was moral valence, include good,  
937        neutral and bad; the third variable was the matching between shape and label: match  
938        vs. nonmatch.

939        ***Stimuli.***

940        The stimuli used in study 3a share the same parameters with experiment 1 & 2. The  
941        differences was that we used six shapes: triangle, square, circle, trapezoid, diamond,  
942        regular pentagon, and six labels: good self, neutral self, bad self, good person, bad person,  
943        and neutral person. To match the concreteness of the label, we asked participant to chosen  
944        an unfamiliar name of their own gender to be the stranger.

945        ***Procedure.***

946        After being fully explained and signed the informed consent, participants were  
947        instructed to chose a name that can represent a stranger with same gender as the  
948        participant themselves, from a common Chinese name pool. Before experiment, the  
949        experimenter explained the meaning of each label to participants. For example, the "good  
950        self" mean the morally good side of themselves, them could imagine the moment when they  
951        do something's morally applauded, "bad self" means the morally bad side of themselves,  
952        they could also imagine the moment when they doing something morally wrong, and

953 “neutral self” means the aspect of self that does not relate to morality, they could imagine  
954 the moment when they doing something irrelevant to morality. In the same sense, the  
955 “good other,” “bad other,” and “neutral other” means the three different aspects of the  
956 stranger, whose name was chosen before the experiment. Then, the experiment proceeded  
957 as study 1a. Each participant finished 6 blocks, each have 120 trials. The sequence of trials  
958 was pseudo-randomized so that there are 10 matched trials for each condition and 10  
959 non-matched trials for each condition (good self, neutral self, bad self, good other, neutral  
960 other, bad other) for each block.

961 ***Data Analysis.***

962 Data analysis followed strategies described in the general method section. Reaction  
963 times and  $d$  prime data were analyzed as in study 1 and study 2, except that one more  
964 within-subject variable (i.e., self-relevance) was included in the analysis.

965 **Results.**

966 ***NHST.***

967 Figure 3 shows  $d$  prime and reaction times of experiment 3a. Less than 5% correct  
968 trials with less than 200ms reaction times were excluded.

969  $d$  prime.

970 There was evidence for the main effect of valence,  $F(1.89, 64.37) = 11.09$ ,  
971  $MSE = 0.23$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .039$ , and main effect of self-relevance,  $F(1, 34) = 3.22$ ,  
972  $MSE = 0.54$ ,  $p = .082$ ,  $\hat{\eta}_G^2 = .015$ , as well as the interaction,  $F(1.79, 60.79) = 3.39$ ,  
973  $MSE = 0.43$ ,  $p = .045$ ,  $\hat{\eta}_G^2 = .022$ .

974 We then conducted separated ANOVA for self-referential and other-referential trials.  
975 The valence effect was shown for the self-referential conditions,  $F(1.65, 56.25) = 13.98$ ,  
976  $MSE = 0.31$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .119$ . Post-hoc test revealed that the Good-Self condition  
977 ( $1.97 \pm 0.14$ ) was with greater  $d$  prime than Netural condition ( $1.41 \pm 0.12$ ,  $t(34) = 4.505$ ,

978  $p = 0.0002$ ), and Bad-self condition ( $1.43 \pm 0.102$ ),  $t(34) = 3.856$ ,  $p = 0.0014$ . There was  
 979 difference between neutral and bad condition,  $t(34) = -0.238$ ,  $p = 0.9694$ . However, no  
 980 effect of valence was found for the other-referential condition  $F(1.98, 67.36) = 0.38$ ,  
 981  $MSE = 0.35$ ,  $p = .681$ ,  $\hat{\eta}_G^2 = .004$ .

982 *Reaction time.*

983 We found interaction between Matchness and Valence ( $F(1.98, 67.44) = 26.29$ ,  
 984  $MSE = 730.09$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .025$ ) and then analyzed the matched trials and nonmatch  
 985 trials separately, as in previous experiments.

986 For the match trials, we found that the interaction between identity and valence,  
 987  $F(1.72, 58.61) = 3.89$ ,  $MSE = 2,750.19$ ,  $p = .032$ ,  $\hat{\eta}_G^2 = .019$ , as well as the main effect of  
 988 valence  $F(1.98, 67.34) = 35.76$ ,  $MSE = 1,127.25$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .079$ , but not the effect of  
 989 identity  $F(1, 34) = 0.20$ ,  $MSE = 3,507.14$ ,  $p = .660$ ,  $\hat{\eta}_G^2 = .001$ . As for the  $d$  prime, we  
 990 separated analyzed the self-referential and other-referential trials. For the Self-referential  
 991 trials, we found the main effect of valence,  $F(1.80, 61.09) = 30.39$ ,  $MSE = 1,584.53$ ,  
 992  $p < .001$ ,  $\hat{\eta}_G^2 = .159$ ; for the other-referential trials, the effect of valence is weaker,  
 993  $F(1.86, 63.08) = 2.85$ ,  $MSE = 2,224.30$ ,  $p = .069$ ,  $\hat{\eta}_G^2 = .024$ . We then focused on the self  
 994 conditions: the good-self condition ( $713 \pm 12$ ) is faster than neutral- ( $776 \pm 11.8$ ),  $t(34) =$   
 995  $-7.396$ ,  $p < .0001$ , and bad-self ( $772 \pm 10.1$ ) conditions,  $t(34) = -5.66$ ,  $p < .0001$ . But  
 996 there is not difference between neutral- and bad-self conditions,  $t(34) = 0.481$ ,  $p = 0.881$ .

997 For the nonmatch trials, we didn't found any strong effect: identity,  $F(1, 34) = 3.43$ ,  
 998  $MSE = 660.02$ ,  $p = .073$ ,  $\hat{\eta}_G^2 = .004$ , valence  $F(1.89, 64.33) = 0.40$ ,  $MSE = 444.10$ ,  
 999  $p = .661$ ,  $\hat{\eta}_G^2 = .001$ , or interaction between the two  $F(1.94, 66.02) = 2.42$ ,  $MSE = 817.35$ ,  
 1000  $p = .099$ ,  $\hat{\eta}_G^2 = .007$ .

1001 **BGLM.**

1002 *Signal detection theory analysis of accuracy.*

1003 We found that the  $d$  prime is greater when shapes were associated with good self

1004 condition than with neutral self or bad self, but shapes associated with bad self and neutral  
1005 self didn't show differences. Comparing the self vs other under three condition revealed  
1006 that shapes associated with good self is greater than with good other, but with a weak  
1007 evidence. In contrast, for both neutral and bad valence condition, shapes associated with  
1008 other had greater  $d$  prime than with self.

1009 *Reaction time.*

1010 In reaction times, we found that same trends in the match trials as in the RT: while  
1011 the shapes associated with good self was greater than with good other (log mean diff =  
1012 -0.02858, 95%HPD[-0.070898, 0.0154]), the direction is reversed for neutral and negative  
1013 condition. see Figure 4

1014 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu,  
1015 Lan, Macrae, & Sui, 2020). We estimated separate drift rate ( $v$ ), non-decision time ( $T_0$ ),  
1016 and boundary separation ( $a$ ) for each condition. We found that the shapes tagged with  
1017 good person has higher drift rate and higher boundary separation than shapes tagged with  
1018 both neutral and bad person. Also, the shapes tagged with neutral person has a higher drift  
1019 rate than shapes tagged with bad person, but not for the boundary separation. Finally, we  
1020 found that shapes tagged with bad person had longer non-decision time (see figure 5)).

1021 **Experiment 3b**

1022 In study 3a, participants had to remember 6 pairs of association, which cause high  
1023 cognitive load during the whole experiment. To eliminate the influence of cognitive load, we  
1024 conducted study 3b, in which participant learn three aspect of self and stranger separately  
1025 in to consecutive task. We hypothesize that we will replicate the pattern of study 3a, i.e.,  
1026 the effect of moral valence only occurs for self-relevant conditions. ### Method

1027 **Participants.**

1028 Study 3b were finished in 2017, at that time we have calculated that the effect size

1029 (Cohen's  $d$ ) of good-person (or good-self) vs. bad-person (or bad-other) was between  $0.47 \sim$   
1030 0.53, based on the data from Tsinghua community in study 1a, 1b, 2, 3a, 4a, and 4b. Based  
1031 on this effect size, we estimated that 54 participants would allow us to detect the effect  
1032 size of Cohen's  $= 0.5$  with 95% power and alpha = 0.05, using G\*power 3.192 (Faul,  
1033 Erdfelder, Buchner, & Lang, 2009). Therefore, we planned to stop after we arrived this  
1034 number. During the data collected at Wenzhou University, 61 participants (45 females; 19  
1035 to 25 years of age, age =  $20.42 \pm 1.77$ ) came to the testing room and we tested all of them  
1036 during a single day. All participants were right-handed, and all had normal or  
1037 corrected-to-normal vision. Informed consent was obtained from all participants prior to  
1038 the experiment according to procedures approved by a local ethics committee. 4  
1039 participants' data were excluded from analysis because their overall accuracy was lower  
1040 than 60%, 1 more participant was excluded because of zero hit rate for one condition,  
1041 leaving 56 participants (43 females; 19 to 25 years old, age =  $20.27 \pm 1.60$ ).

1042 ***Design.***

1043 Study 3b has the same experimental design as 3a, with a  $2 \times 3 \times 2$  within-subject  
1044 design. The first variable was self-relevance, include two levels: self-relevant  
1045 vs. stranger-relevant; the second variable was moral valence, include good, neutral and bad;  
1046 the third variable was the matching between shape and label: match vs. mismatch.  
1047 Stimuli. The stimuli used in study 3b share the same parameters with experiment 3a. 6  
1048 shapes were included (triangle, square, circle, trapezoid, diamond, regular pentagon), as  
1049 well as 6 labels, but the labels changed to "good self," "neutral self," "bad self," "good  
1050 him/her," "bad him/her", "neutral him/her," the stranger's label is consistent with  
1051 participants' gender. Same as study 3a, we asked participant to chosen an unfamiliar name  
1052 of their own gender to be the stranger before showing them the relationship. Note, because  
1053 of implementing error, the personal distance data did not collect for this experiment.

1054 ***Stimuli.***

1055 The stimuli used in study 3b is the same as in experiment 3a.

1056 ***Procedure.***

1057 In this experiment, participants finished two matching tasks, i.e., self-matching task,  
1058 and other-matching task. In the self-matching task, participants first associate the three  
1059 aspects of self to three different shapes, and then perform the matching task. In the  
1060 other-matching task, participants first associate the three aspects of the stranger to three  
1061 different shapes, and then perform the matching task. The order of self-task and other-task  
1062 are counter-balanced among participants. Different from experiment 3a, after presenting  
1063 the stimuli pair for 100ms, participant has 1900 ms to response, and they feedback with  
1064 both accuracy and reaction time. As in study 3a, before each task, the instruction showed  
1065 the meaning of each label to participants. The self-matching task and other-matching task  
1066 were randomized between participants. Each participant finished 6 blocks, each have 120  
1067 trials.

1068 ***Data Analysis.***

1069 Same as experiment 3a.

1070 ***Results.***

1071 ***NHST.***

1072 Figure 6 shows  $d$  prime and reaction times of experiment 3b. Less than 5% correct  
1073 trials with less than 200ms reaction times were excluded.

1074  $d$  prime.

1075 There was no evidence for the main effect of valence,  $F(1.92, 105.43) = 1.90$ ,  
1076  $MSE = 0.33$ ,  $p = .157$ ,  $\hat{\eta}_G^2 = .005$ , but we found a main effect of self-relevance,  
1077  $F(1, 55) = 4.65$ ,  $MSE = 0.89$ ,  $p = .035$ ,  $\hat{\eta}_G^2 = .017$ , as well as the interaction,  
1078  $F(1.90, 104.36) = 5.58$ ,  $MSE = 0.26$ ,  $p = .006$ ,  $\hat{\eta}_G^2 = .011$ .

1079 We then conducted separated ANOVA for self-referential and other-referential trials.

1080 The valence effect was shown for the self-referential conditions,  $F(1.75, 96.42) = 6.73$ ,  
 1081  $MSE = 0.30$ ,  $p = .003$ ,  $\hat{\eta}_G^2 = .037$ . Post-hoc test revealed that the Good-Self condition  
 1082 ( $2.15 \pm 0.12$ ) was with greater  $d$  prime than Neutral condition ( $1.83 \pm 0.12$ ,  $t(34) = 3.36$ ,  
 1083  $p = 0.0031$ ), and Bad-self condition ( $1.87 \pm 0.12$ ),  $t(34) = 2.955$ ,  $p = 0.01$ . There was  
 1084 difference between neutral and bad condition,  $t(34) = -0.039$ ,  $p = 0.914$ . However, no effect  
 1085 of valence was found for the other-referential condition  $F(1.93, 105.97) = 0.61$ ,  
 1086  $MSE = 0.31$ ,  $p = .539$ ,  $\hat{\eta}_G^2 = .002$ .

1087 *Reaction time.*

1088 We found interaction between Matchness and Valence ( $F(1.86, 102.47) = 15.44$ ,  
 1089  $MSE = 3,112.78$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .006$ ) and then analyzed the matched trials and  
 1090 nonmatch trials separately, as in previous experiments.

1091 For the match trials, we found that the interaction between identity and valence,  
 1092  $F(1.67, 92.11) = 6.14$ ,  $MSE = 6,472.48$ ,  $p = .005$ ,  $\hat{\eta}_G^2 = .009$ , as well as the main effect of  
 1093 valence  $F(1.88, 103.65) = 24.25$ ,  $MSE = 5,994.25$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .038$ , but not the effect  
 1094 of identity  $F(1, 55) = 48.49$ ,  $MSE = 25,892.59$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .153$ . As for the  $d$  prime,  
 1095 we separated analyzed the self-referential and other-referential trials. For the  
 1096 Self-referential trials, we found the main effect of valence,  $F(1.66, 91.38) = 23.98$ ,  
 1097  $MSE = 6,965.61$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .100$ ; for the other-referential trials, the effect of valence  
 1098 is weaker,  $F(1.89, 103.94) = 5.96$ ,  $MSE = 5,589.90$ ,  $p = .004$ ,  $\hat{\eta}_G^2 = .014$ . We then focused  
 1099 on the self conditions: the good-self condition ( $713 \pm 12$ ) is faster than neutral- ( $776 \pm$   
 1100  $11.8$ ),  $t(34) = -7.396$ ,  $p < .0001$ , and bad-self ( $772 \pm 10.1$ ) conditions,  $t(34) = -5.66$ ,  $p <$   
 1101  $.0001$ . But there is not difference between neutral- and bad-self conditions,  $t(34) = 0.481$ ,  $p$   
 1102  $= 0.881$ .

1103 For the nonmatch trials, we didn't found any strong effect: identity,  $F(1, 55) = 10.31$ ,  
 1104  $MSE = 24,590.52$ ,  $p = .002$ ,  $\hat{\eta}_G^2 = .035$ , valence  $F(1.98, 108.63) = 20.57$ ,  $MSE = 2,847.51$ ,  
 1105  $p < .001$ ,  $\hat{\eta}_G^2 = .016$ , or interaction between the two  $F(1.93, 106.25) = 35.51$ ,

<sub>1106</sub>  $MSE = 1,939.88$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .019$ .

<sub>1107</sub> **BGLM.**

<sub>1108</sub> *Signal detection theory analysis of accuracy.*

<sub>1109</sub> We found that the  $d$  prime is greater when shapes were associated with good self  
<sub>1110</sub> condition than with neutral self or bad self, but shapes associated with bad self and neutral  
<sub>1111</sub> self didn't show differences. comparing the self vs other under three condition revealed that  
<sub>1112</sub> shapes associated with good self is greater than with good other, but with a weak evidence.  
<sub>1113</sub> In contrast, for both neutral and bad valence condition, shapes associated with other had  
<sub>1114</sub> greater  $d$  prime than with self.

<sub>1115</sub> *Reaction time.*

<sub>1116</sub> In reaction times, we found that same trends in the match trials as in the RT: while  
<sub>1117</sub> the shapes associated with good self was greater than with good other (log mean diff =  
<sub>1118</sub> -0.02858, 95%HPD[-0.070898, 0.0154]), the direction is reversed for neutral and negative  
<sub>1119</sub> condition. see Figure 7

<sub>1120</sub> **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu,  
<sub>1121</sub> Lan, Macrae, & Sui, 2020). We estimated separate drift rate ( $v$ ), non-decision time ( $T_0$ ),  
<sub>1122</sub> and boundary separation ( $a$ ) for each condition. We found that, similar to experiment 3a,  
<sub>1123</sub> the shapes tagged with good person has higher drift rate and higher boundary separation  
<sub>1124</sub> than shapes tagged with both neutral and bad person, but only for the self-referential  
<sub>1125</sub> condition. Also, the shapes tagged with neutral person has a higher drift rate than shapes  
<sub>1126</sub> tagged with bad person, but not for the boundary separation, and this effect also exist only  
<sub>1127</sub> for the self-referential condition.

<sub>1128</sub> Interestingly, we found that in both self-referential and other-referential conditions,  
<sub>1129</sub> the shapes associated bad valence have higher drift rate and higher boundary separation.  
<sub>1130</sub> which might suggest that the shape associated with bad stimuli might be prioritized in the  
<sub>1131</sub> non-match trials (see figure 8)).

1132 **Experiment 6b**

1133       Experiment 6b was conducted to study the neural correlates of the prioritization  
1134       effect of positive self, i.e., the neural underlying of the behavioral effect found in  
1135       experiment 3a. However, as in experiment 6a, the procedure of this experiment was  
1136       modified to adopted to ERP experiment.

1137       **Method.**

1138       ***Participants.***

1139       23 college students (8 female, age =  $22.86 \pm 2.47$ ) participated the current study, all  
1140       of them were recruited from Tsinghua University in 2016. Informed consent was obtained  
1141       from all participants prior to the experiment according to procedures approved by a local  
1142       ethics committee. For day 1's data, 1 participant was excluded from the current analysis  
1143       because of lower than 60% overall accuracy, remaining 22 participants (8 female, age =  
1144        $22.76 \pm 2.49$ ). For day 2's data, one participant dropped out, leaving 22 participants (9  
1145       female, age =  $23.05 \pm 2.46$ ), all of them has overall accuracy higher than 60%.

1146       ***Design.***

1147       The experimental design of this experiment is same as experiment 3: a  $2 \times 3 \times 2$   
1148       within-subject design with self-relevance (self-relevant vs. other-relevant), moral valence  
1149       (good, neutral, and bad) and matchness between shape and label (match vs. mismatch) as  
1150       within-subject variables.

1151       ***Stimuli.***

1152       As in experiment 3a, 6 shapes were included (triangle, square, circle, trapezoid,  
1153       diamond, regular pentagon), as well as 6 labels (good self, neutral self, bad self, good  
1154       person, bad person, neutral person). To match the concreteness of the label, we asked  
1155       participant to chosen an unfamiliar name of their own gender to be the stranger.

1156        ***Procedure.***

1157        The procedure was similar to Experiment 2 and 6a. Subjects first learned the  
1158        associations between labels and shapes and then completed a shape-label matching task. In  
1159        each trial of the matching task, a fixation were first presented for 500 ms, followed by a 50  
1160        ms label; then, a scrambled picture presented 900 ms. After the backward mask, the shape  
1161        were presented on a noisy background for 50ms. Participant have to response in 1000ms  
1162        after the presentation of the shape, and finally, a feedback screen was presented for 500 ms.  
1163        The inter-trial interval (ITI) were randomly varied at the range of 1000 ~ 1400 ms.

1164        All the stimuli were presented on a gray background (RGB: 127, 127, 127). E-primed  
1165        2.0 was used to present stimuli and collect behavioral results. Data were collected and  
1166        analyzed when accuracy performance in total reached 60%.

1167        Because learning 6 associations was more difficult than 3 associations and participant  
1168        might have low accuracy (see experiment 3a), the current study had extended to a two-day  
1169        paradigm to maximizing the accurate trials that can be used in EEG data. At the first day,  
1170        participants learnt the associations and finished 9 blocks of the matching task, each had  
1171        120 trials, without EEG recording. That is, each condition has 90 trials.

1172        Participants came back to lab at the second day and finish the same task again, with  
1173        EEG recorded. Before the EEG experiment, each participant finished a practice session  
1174        again, if their accuracy is equal or higher than 85%, they start the experiment (one  
1175        participant used lower threshold 75%). Each participant finished 18 blocks, each has 90  
1176        trials. One participant finished additional 6 blocks because of high error rate at the  
1177        beginning, another two participant finished addition 3 blocks because of the technique  
1178        failure in recording the EEG data. To increase the number of trials that can be used for  
1179        EEG data analysis, matched trials has twice number as mismatched trials, therefore, for  
1180        matched trials each participants finished 180 trials for each condition, for mismatched  
1181        trials, each conditions has 90 trials.

1182     ***Data Analysis.***

1183     Same as experiment 3a.

1184     **Results of Day 1.**

1185     ***NHST.***

1186     Figure 9 shows  $d$  prime and reaction times of experiment 3b. Less than 5% correct  
 1187     trials with less than 200ms reaction times were excluded.

1188      $d$  prime.

1189     There was no evidence for the main effect of valence,  $F(1.91, 40.20) = 11.98$ ,

1190      $MSE = 0.15$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .040$ , but we found a main effect of self-relevance,

1191      $F(1, 21) = 1.21$ ,  $MSE = 0.20$ ,  $p = .284$ ,  $\hat{\eta}_G^2 = .003$ , as well as the interaction,

1192      $F(1.28, 26.90) = 12.88$ ,  $MSE = 0.21$ ,  $p = .001$ ,  $\hat{\eta}_G^2 = .041$ .

1193     We then conducted separated ANOVA for self-referential and other-referential trials.

1194     The valence effect was shown for the self-referential conditions,  $F(1.73, 36.42) = 29.31$ ,

1195      $MSE = 0.14$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .147$ . Post-hoc test revealed that the Good-Self condition

1196      $(2.15 \pm 0.12)$  was with greater  $d$  prime than Neutral condition  $(1.83 \pm 0.12$ ,  $t(34) = 3.36$ ,

1197      $p = 0.0031$ ), and Bad-self condition  $(1.87 \pm 0.12)$ ,  $t(34) = 2.955$ ,  $p = 0.01$ . There was

1198     difference between neutral and bad condition,  $t(34) = -0.039$ ,  $p = 0.914$ . However, no effect

1199     of valence was found for the other-referential condition  $F(1.75, 36.72) = 0.00$ ,  $MSE = 0.18$ ,

1200      $p = .999$ ,  $\hat{\eta}_G^2 = .000$ .

1201     *Reaction time.*

1202     We found interaction between Matchness and Valence ( $F(1.79, 37.63) = 4.07$ ,

1203      $MSE = 704.90$ ,  $p = .029$ ,  $\hat{\eta}_G^2 = .003$ ) and then analyzed the matched trials and nonmatch

1204     trials separately, as in previous experiments.

1205     For the match trials, we found that the interaction between identity and valence,

1206      $F(1.72, 36.16) = 4.55$ ,  $MSE = 1,560.90$ ,  $p = .022$ ,  $\hat{\eta}_G^2 = .015$ , as well as the main effect of

valence  $F(1.93, 40.55) = 9.83$ ,  $MSE = 1,951.84$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .044$ , but not the effect of identity  $F(1, 21) = 4.87$ ,  $MSE = 2,032.05$ ,  $p = .039$ ,  $\hat{\eta}_G^2 = .012$ . As for the  $d$  prime, we separated analyzed the self-referential and other-referential trials. For the Self-referential trials, we found the main effect of valence,  $F(1.92, 40.38) = 14.48$ ,  $MSE = 1,647.20$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .112$ ; for the other-referential trials, the effect of valence is weaker,  $F(1.79, 37.50) = 1.04$ ,  $MSE = 1,842.07$ ,  $p = .356$ ,  $\hat{\eta}_G^2 = .008$ . We then focused on the self conditions: the good-self condition ( $713 \pm 12$ ) is faster than neutral- ( $776 \pm 11.8$ ),  $t(34) = -7.396$ ,  $p < .0001$ , and bad-self ( $772 \pm 10.1$ ) conditions,  $t(34) = -5.66$ ,  $p < .0001$ . But there is not difference between neutral- and bad-self conditions,  $t(34) = 0.481$ ,  $p = 0.881$ .

For the nonmatch trials, we didn't found any strong effect: identity,  $F(1, 21) = 2.76$ ,  $MSE = 1,718.93$ ,  $p = .112$ ,  $\hat{\eta}_G^2 = .006$ , valence  $F(1.61, 33.77) = 3.81$ ,  $MSE = 1,532.21$ ,  $p = .041$ ,  $\hat{\eta}_G^2 = .012$ , or interaction between the two  $F(1.90, 39.97) = 2.23$ ,  $MSE = 720.80$ ,  $p = .123$ ,  $\hat{\eta}_G^2 = .004$ .

## 1220 **BGLM.**

1221 *Signal detection theory analysis of accuracy.*

1222 We found that the  $d$  prime is greater when shapes were associated with good self  
 1223 condition than with neutral self or bad self, but shapes associated with bad self and neutral  
 1224 self didn't show differences. comparing the self vs other under three condition revealed that  
 1225 shapes associated with good self is greater than with good other, but with a weak evidence.  
 1226 In contrast, for both neutral and bad valence condition, shapes associated with other had  
 1227 greater  $d$  prime than with self.

1228 *Reaction time.*

1229 In reaction times, we found that same trends in the match trials as in the RT: while  
 1230 the shapes associated with good self was greater than with good other (log mean diff =  
 1231 -0.02858, 95%HPD[-0.070898, 0.0154]), the direction is reversed for neutral and negative  
 1232 condition. see Figure 10

1233       **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu,

1234 Lan, Macrae, & Sui, 2020). We estimated separate drift rate ( $v$ ), non-decision time ( $T_0$ ),

1235 and boundary separation ( $a$ ) for each condition. We found that, similar to experiment 3a,

1236 the shapes tagged with good person has higher drift rate and higher boundary separation

1237 than shapes tagged with both neutral and bad person, but only for the self-referential

1238 condition. Also, the shapes tagged with neutral person has a higher drift rate than shapes

1239 tagged with bad person, but not for the boundary separation, and this effect also exist only

1240 for the self-referential condition.

1241       Interestingly, we found that in both self-referential and other-referential conditions,

1242 the shapes associated bad valence have higher drift rate and higher boundary separation.

1243 which might suggest that the shape associated with bad stimuli might be prioritized in the

1244 non-match trials (see figure 11).

### 1245       **Part 3: Implicit binding between valence and identity**

1246       In this part, we reported two studies in which the moral valence or the self-referential

1247 processing is not task-relevant. We are interested in testing whether the task-relevance will

1248 eliminate the effect observed in previous experiment.

### 1249       **Experiment 4a: Morality as task-irrelevant variable**

1250       In part two (experiment 3a and 3b), participants learned the association between self

1251 and moral valence directly. In Experiment 4a, we examined whether the interaction

1252 between moral valence and identity occur even when one of the variable was irrelevant to

1253 the task. In experiment 4a, participants learnt associations between shapes and self/other

1254 labels, then made perceptual match judgments only about the self or other conditions

1255 labels and shapes (cf. Sui, He, and Humphreys (2012)). However, we presented labels of

1256 different moral valence in the shapes, which means that the moral valence factor become

1257 task irrelevant. If the binding between moral good and self is intrinsic and automatic, then  
1258 we will observe that facilitating effect of moral good for self conditions, but not for other  
1259 conditions.

1260 **Method.**

1261 ***Participants.***

1262 64 participants (37 female, age =  $19.70 \pm 1.22$ ) participated the current study, 32 of  
1263 them were from Tsinghua University in 2015, 32 were from Wenzhou University  
1264 participated in 2017. All participants were right-handed, and all had normal or  
1265 corrected-to-normal vision. Informed consent was obtained from all participants prior to  
1266 the experiment according to procedures approved by a local ethics committee. The data  
1267 from 5 participants from Wenzhou site were excluded from analysis because their accuracy  
1268 was close to chance ( $< 0.6$ ). The results for the remaining 59 participants (33 female, age  
1269 =  $19.78 \pm 1.20$ ) were analyzed and reported.

1270 ***Design.***

1271 As in Experiment 3, a  $2 \times 3 \times 2$  within-subject design was used. The first variable was  
1272 self-relevance (self and stranger associations); the second variable was moral valence (good,  
1273 neutral and bad associations); the third variable was the matching between shape and label  
1274 (matching vs. non-match for the personal association). However, in this the task,  
1275 participants only learn the association between two geometric shapes and two labels (self  
1276 and other), i.e., only self-relevance were related to the task. The moral valence  
1277 manipulation was achieved by embedding the personal label of the labels in the geometric  
1278 shapes, see below. For simplicity, the trials where shapes where paired with self and with a  
1279 word of “good person” inside were shorted as good-self condition, similarly, the trials where  
1280 shapes paired with the self and with a word of “bad person” inside were shorted as bad-self  
1281 condition. Hence, we also have six conditions: good-self, neutral-self, bad-self, good-other,  
1282 neutral-other, and bad-other.

***Stimuli.***

2 shapes were included (circle, square) and each appeared above a central fixation cross with the personal label appearing below. However, the shapes were not empty but with a two-Chinese-character word in the middle, the word was one of three labels with different moral valence: “good person,” “bad person” and “neutral person.” Before the experiment, participants learned the self/other association, and were informed to only response to the association between shapes’ configures and the labels below the fixation, but ignore the words within shapes. Besides the behavioral experiments, participants from Tsinghua community also finished questionnaires as Experiments 3, and participants from Wenzhou community finished a series of questionnaire as the other experiment finished in Wenzhou.

***Procedure.***

The procedure was similar to Experiment 1. There were 6 blocks of trial, each with 120 trials for 2017 data. Due to procedure error, the data collected in 2015 in Tsinghua community only have 60 trials for each block, i.e., 30 trials per condition.

As in study 3a, before each task, the instruction showed the meaning of each label to participants. The self-matching task and other-matching task were randomized between participants. Each participant finished 6 blocks, each have 120 trials.

***Data Analysis.***

Same as experiment 3a.

***Results.******NHST.***

Figure 12 shows  $d$  prime and reaction times of experiment 3a. Less than 5% correct trials with less than 200ms reaction times were excluded.

$d$  prime.

1308 There was no evidence for the main effect of valence,  $F(1.93, 111.66) = 0.53$ ,  
 1309  $MSE = 0.12, p = .581, \hat{\eta}_G^2 = .000$ , but we found a main effect of self-relevance,  
 1310  $F(1, 58) = 121.04, MSE = 0.48, p < .001, \hat{\eta}_G^2 = .189$ , as well as the interaction,  
 1311  $F(1.99, 115.20) = 4.12, MSE = 0.14, p = .019, \hat{\eta}_G^2 = .004$ .

1312 We then conducted separated ANOVA for self-referential and other-referential trials.  
 1313 The valence effect was shown for the self-referential conditions,  $F(1.95, 112.92) = 3.01$ ,  
 1314  $MSE = 0.15, p = .055, \hat{\eta}_G^2 = .008$ . Post-hoc test revealed that the Good-Self condition  
 1315 ( $2.15 \pm 0.12$ ) was with greater  $d$  prime than Neutral condition ( $1.83 \pm 0.12, t(34) = 3.36$ ,  
 1316  $p = 0.0031$ ), and Bad-self condition ( $1.87 \pm 0.12$ ),  $t(34) = 2.955, p = 0.01$ . There was  
 1317 difference between neutral and bad condition,  $t(34) = -0.039, p = 0.914$ . However, no effect  
 1318 of valence was found for the other-referential condition  $F(1.98, 114.61) = 1.75$ ,  
 1319  $MSE = 0.10, p = .179, \hat{\eta}_G^2 = .003$ .

1320 *Reaction time.*

1321 We found interaction between Matchness and Valence ( $F(1.94, 112.64) = 0.84$ ,  
 1322  $MSE = 465.35, p = .432, \hat{\eta}_G^2 = .000$ ) and then analyzed the matched trials and nonmatch  
 1323 trials separately, as in previous experiments.

1324 For the match trials, we found that the interaction between identity and valence,  
 1325  $F(1.90, 110.18) = 4.41, MSE = 465.91, p = .016, \hat{\eta}_G^2 = .003$ , as well as the main effect of  
 1326 valence  $F(1.98, 114.82) = 0.94, MSE = 606.30, p = .392, \hat{\eta}_G^2 = .001$ , but not the effect of  
 1327 identity  $F(1, 58) = 124.15, MSE = 4,037.53, p < .001, \hat{\eta}_G^2 = .257$ . As for the  $d$  prime, we  
 1328 separated analyzed the self-referential and other-referential trials. For the Self-referential  
 1329 trials, we found the main effect of valence,  $F(1.97, 114.32) = 6.29, MSE = 367.25$ ,  
 1330  $p = .003, \hat{\eta}_G^2 = .006$ ; for the other-referential trials, the effect of valence is weaker,  
 1331  $F(1.95, 112.89) = 0.35, MSE = 699.50, p = .699, \hat{\eta}_G^2 = .001$ . We then focused on the self  
 1332 conditions: the good-self condition ( $713 \pm 12$ ) is faster than neutral- ( $776 \pm 11.8$ ),  $t(34) =$   
 1333  $-7.396, p < .0001$ , and bad-self ( $772 \pm 10.1$ ) conditions,  $t(34) = -5.66, p < .0001$ . But

<sup>1334</sup> there is not difference between neutral- and bad-self conditions,  $t(34) = 0.481$ ,  $p = 0.881$ .

<sup>1335</sup> For the nonmatch trials, we didn't found any strong effect: identity,  $F(1, 58) = 0.16$ ,

<sup>1336</sup>  $MSE = 1,547.37$ ,  $p = .692$ ,  $\hat{\eta}_G^2 = .000$ , valence  $F(1.96, 113.52) = 0.68$ ,  $MSE = 390.26$ ,

<sup>1337</sup>  $p = .508$ ,  $\hat{\eta}_G^2 = .000$ , or interaction between the two  $F(1.90, 110.27) = 0.04$ ,

<sup>1338</sup>  $MSE = 585.80$ ,  $p = .953$ ,  $\hat{\eta}_G^2 = .000$ .

<sup>1339</sup> **BGLM.**

<sup>1340</sup> *Signal detection theory analysis of accuracy.*

<sup>1341</sup> We found that the  $d$  prime is greater when shapes were associated with good self  
<sup>1342</sup> condition than with neutral self or bad self, but shapes associated with bad self and neutral  
<sup>1343</sup> self didn't show differences. comparing the self vs other under three condition revealed that  
<sup>1344</sup> shapes associated with good self is greater than with good other, but with a weak evidence.

<sup>1345</sup> In contrast, for both neutral and bad valence condition, shapes associated with other had  
<sup>1346</sup> greater  $d$  prime than with self.

<sup>1347</sup> *Reaction time.*

<sup>1348</sup> In reaction times, we found that same trends in the match trials as in the RT: while  
<sup>1349</sup> the shapes associated with good self was greater than with good other ( $\log \text{mean diff} =$   
<sup>1350</sup>  $-0.02858$ ,  $95\% \text{HPD}[-0.070898, 0.0154]$ ), the direction is reversed for neutral and negative  
<sup>1351</sup> condition. see Figure 13

<sup>1352</sup> **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu,  
<sup>1353</sup> Lan, Macrae, & Sui, 2020). We estimated separate drift rate ( $v$ ), non-decision time ( $T_0$ ),  
<sup>1354</sup> and boundary separation ( $a$ ) for each condition. We found that the shapes tagged with  
<sup>1355</sup> good person has higher drift rate and higher boundary separation than shapes tagged with  
<sup>1356</sup> both neutral and bad person. Also, the shapes tagged with neutral person has a higher drift  
<sup>1357</sup> rate than shapes tagged with bad person, but not for the boundary separation. Finally, we  
<sup>1358</sup> found that shapes tagged with bad person had longer non-decision time (see figure 14)).

**1359 Experiment 4b: Morality as task-irrelevant variable**

1360 In study 4b, we changed the role of valence and identity in task. In this experiment,  
1361 participants learn the association between moral valence and the made perceptual match  
1362 judgments to associations between different moral valence and shapes as in study 1-3.  
1363 Different from experiment 1 ~ 3, we made put the labels of “self/other” in the shapes so  
1364 that identity served as an task irrelevant variable. As in experiment 4b, we also  
1365 hypothesized that the intrinsic binding between morally good and self will enhance the  
1366 performance of good self condition, even identity is irrelevant to the task.

1367 **Method.**

1368 ***Participants.***

1369 53 participants (39 female, age =  $20.57 \pm 1.81$ ) participated the current study, 34 of  
1370 them were from Tsinghua University in 2015, 19 were from Wenzhou University  
1371 participated in 2017. All participants were right-handed, and all had normal or  
1372 corrected-to-normal vision. Informed consent was obtained from all participants prior to  
1373 the experiment according to procedures approved by a local ethics committee. The data  
1374 from 8 participants from Wenzhou site were excluded from analysis because their accuracy  
1375 was close to chance ( $< 0.6$ ). The results for the remaining 45 participants (33 female, age  
1376 =  $20.78 \pm 1.76$ ) were analyzed and reported.

1377 ***Design.***

1378 As in Experiment 3, a  $2 \times 3 \times 2$  within-subject design was used. The first variable was  
1379 self-relevance (self and stranger associations); the second variable was moral valence (good,  
1380 neutral and bad associations); the third variable was the matching between shape and label  
1381 (matching vs. non-match for the personal association). However, in this the task,  
1382 participants only learn the association between two geometric shapes and two labels (self  
1383 and other), i.e., only self-relevance were related to the task. The moral valence  
1384 manipulation was achieved by embedding the personal label of the labels in the geometric

1385 shapes, see below. For simplicity, the trials where shapes were paired with self and with a  
1386 word of “good person” inside were shorted as good-self condition, similarly, the trials where  
1387 shapes paired with the self and with a word of “bad person” inside were shorted as bad-self  
1388 condition. Hence, we also have six conditions: good-self, neutral-self, bad-self, good-other,  
1389 neutral-other, and bad-other.

1390 ***Stimuli.***

1391 2 shapes were included (circle, square) and each appeared above a central fixation  
1392 cross with the personal label appearing below. However, the shapes were not empty but  
1393 with a two-Chinese-character word in the middle, the word was one of three labels with  
1394 different moral valence: “good person,” “bad person” and “neutral person.” Before the  
1395 experiment, participants learned the self/other association, and were informed to only  
1396 response to the association between shapes’ configures and the labels below the fixation, but  
1397 ignore the words within shapes. Besides the behavioral experiments, participants from  
1398 Tsinghua community also finished questionnaires as Experiments 3, and participants from  
1399 Wenzhou community finished a series of questionnaire as the other experiment finished in  
1400 Wenzhou.

1401 ***Procedure.***

1402 The procedure was similar to Experiment 1. There were 6 blocks of trial, each with  
1403 120 trials for 2017 data. Due to procedure error, the data collected in 2015 in Tsinghua  
1404 community only have 60 trials for each block, i.e., 30 trials per condition.

1405 As in study 3a, before each task, the instruction showed the meaning of each label to  
1406 participants. The self-matching task and other-matching task were randomized between  
1407 participants. Each participant finished 6 blocks, each have 120 trials.

1408 ***Data Analysis.***

1409 Same as experiment 3a.

1410 ***Results.***

<sup>1411</sup> **NHST.**

<sup>1412</sup> Figure 15 shows  $d$  prime and reaction times of experiment 3a. Less than 5% correct  
<sup>1413</sup> trials with less than 200ms reaction times were excluded.

<sup>1414</sup>  $d$  prime.

<sup>1415</sup> There was no evidence for the main effect of valence,  $F(1.59, 69.94) = 2.34$ ,  
<sup>1416</sup>  $MSE = 0.48$ ,  $p = .115$ ,  $\hat{\eta}_G^2 = .010$ , but we found a main effect of self-relevance,  
<sup>1417</sup>  $F(1, 44) = 0.00$ ,  $MSE = 0.08$ ,  $p = .994$ ,  $\hat{\eta}_G^2 = .000$ , as well as the interaction,  
<sup>1418</sup>  $F(1.96, 86.41) = 0.53$ ,  $MSE = 0.10$ ,  $p = .585$ ,  $\hat{\eta}_G^2 = .001$ .

<sup>1419</sup> We then conducted separated ANOVA for self-referential and other-referential trials.

<sup>1420</sup> The valence effect was shown for the self-referential conditions,  $F(1.75, 76.86) = 3.08$ ,  
<sup>1421</sup>  $MSE = 0.25$ ,  $p = .058$ ,  $\hat{\eta}_G^2 = .017$ . Post-hoc test revealed that the Good-Self condition  
<sup>1422</sup> ( $2.15 \pm 0.12$ ) was with greater  $d$  prime than Neutral condition ( $1.83 \pm 0.12$ ,  $t(34) = 3.36$ ,  
<sup>1423</sup>  $p = 0.0031$ ), and Bad-self condition ( $1.87 \pm 0.12$ ),  $t(34) = 2.955$ ,  $p = 0.01$ . There was  
<sup>1424</sup> difference between neutral and bad condition,  $t(34) = -0.039$ ,  $p = 0.914$ . However, no effect  
<sup>1425</sup> of valence was found for the other-referential condition  $F(1.63, 71.50) = 1.07$ ,  $MSE = 0.33$ ,  
<sup>1426</sup>  $p = .336$ ,  $\hat{\eta}_G^2 = .006$ .

<sup>1427</sup> *Reaction time.*

<sup>1428</sup> We found interaction between Matchness and Valence ( $F(1.87, 82.50) = 18.58$ ,  
<sup>1429</sup>  $MSE = 1,291.12$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .023$ ) and then analyzed the matched trials and  
<sup>1430</sup> nonmatch trials separately, as in previous experiments.

<sup>1431</sup> For the match trials, we found that the interaction between identity and valence,  
<sup>1432</sup>  $F(1.86, 81.84) = 5.22$ ,  $MSE = 308.30$ ,  $p = .009$ ,  $\hat{\eta}_G^2 = .003$ , as well as the main effect of  
<sup>1433</sup> valence  $F(1.80, 79.37) = 11.04$ ,  $MSE = 2,937.54$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .059$ , but not the effect of  
<sup>1434</sup> identity  $F(1, 44) = 0.23$ ,  $MSE = 263.26$ ,  $p = .632$ ,  $\hat{\eta}_G^2 = .000$ . As for the  $d$  prime, we  
<sup>1435</sup> separated analyzed the self-referential and other-referential trials. For the Self-referential  
<sup>1436</sup> trials, we found the main effect of valence,  $F(1.74, 76.48) = 13.69$ ,  $MSE = 1,732.08$ ,

<sup>1437</sup>  $p < .001$ ,  $\hat{\eta}_G^2 = .079$ ; for the other-referential trials, the effect of valence is weaker,  
<sup>1438</sup>  $F(1.87, 82.44) = 7.09$ ,  $MSE = 1,527.43$ ,  $p = .002$ ,  $\hat{\eta}_G^2 = .043$ . We then focused on the self  
<sup>1439</sup> conditions: the good-self condition ( $713 \pm 12$ ) is faster than neutral- ( $776 \pm 11.8$ ),  $t(34) =$   
<sup>1440</sup>  $-7.396$ ,  $p < .0001$ , and bad-self ( $772 \pm 10.1$ ) conditions,  $t(34) = -5.66$ ,  $p < .0001$ . But  
<sup>1441</sup> there is not difference between neutral- and bad-self conditions,  $t(34) = 0.481$ ,  $p = 0.881$ .

<sup>1442</sup> For the nonmatch trials, we didn't found any strong effect: identity,  $F(1, 44) = 1.96$ ,  
<sup>1443</sup>  $MSE = 319.47$ ,  $p = .169$ ,  $\hat{\eta}_G^2 = .001$ , valence  $F(1.69, 74.54) = 6.59$ ,  $MSE = 886.19$ ,  
<sup>1444</sup>  $p = .004$ ,  $\hat{\eta}_G^2 = .010$ , or interaction between the two  $F(1.88, 82.57) = 0.31$ ,  $MSE = 316.96$ ,  
<sup>1445</sup>  $p = .718$ ,  $\hat{\eta}_G^2 = .000$ .

<sup>1446</sup> **BGLM.**

<sup>1447</sup> *Signal detection theory analysis of accuracy.*

<sup>1448</sup> We found that the  $d$  prime is greater when shapes were associated with good self  
<sup>1449</sup> condition than with neutral self or bad self, but shapes associated with bad self and neutral  
<sup>1450</sup> self didn't show differences. comparing the self vs other under three condition revealed that  
<sup>1451</sup> shapes associated with good self is greater than with good other, but with a weak evidence.  
<sup>1452</sup> In contrast, for both neutral and bad valence condition, shapes associated with other had  
<sup>1453</sup> greater  $d$  prime than with self.

<sup>1454</sup> *Reaction time.*

<sup>1455</sup> In reaction times, we found that same trends in the match trials as in the RT: while  
<sup>1456</sup> the shapes associated with good self was greater than with good other (log mean diff =  
<sup>1457</sup>  $-0.02858$ , 95%HPD[-0.070898, 0.0154]), the direction is reversed for neutral and negative  
<sup>1458</sup> condition. see Figure 16

<sup>1459</sup> **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu,  
<sup>1460</sup> Lan, Macrae, & Sui, 2020). We estimated separate drift rate ( $v$ ), non-decision time ( $T_0$ ),  
<sup>1461</sup> and boundary separation ( $a$ ) for each condition. We found that the shapes tagged with  
<sup>1462</sup> good person has higher drift rate and higher boundary separation than shapes tagged with

<sub>1463</sub> both neutral and bad person. Also, the shapes tagged with neutral person has a higher drift  
<sub>1464</sub> rate than shapes tagged with bad person, but not for the boundary separation. Finally, we  
<sub>1465</sub> found that shapes tagged with bad person had longer non-decision time (see figure 17)).

<sub>1466</sub>

## Results

<sub>1467</sub> **Effect of moral valence**

<sub>1468</sub> In this part, we synthesized results from experiment 1a, 1b, 1c, 2, 5 and 6a. Data  
<sub>1469</sub> from 192 participants were included in these analyses. We found differences between  
<sub>1470</sub> positive and negative conditions on RT was Cohen's  $d = -0.58 \pm 0.06$ , 95% CI [-0.70 -0.47];  
<sub>1471</sub> on  $d'$  was Cohen's  $d = 0.24 \pm 0.05$ , 95% CI [0.15 0.34]. The effect was also observed  
<sub>1472</sub> between positive and neutral condition, RT: Cohen's  $d = -0.44 \pm 0.10$ , 95% CI [-0.63  
<sub>1473</sub> -0.25];  $d'$ : Cohen's  $d = 0.31 \pm 0.07$ , 95% CI [0.16 0.45]. And the difference between neutral  
<sub>1474</sub> and bad conditions are not significant, RT: Cohen's  $d = 0.15 \pm 0.07$ , 95% CI [0.00 0.30];  
<sub>1475</sub>  $d'$ : Cohen's  $d = 0.07 \pm 0.07$ , 95% CI [-0.08 0.21]. See Figure 18 left panel.

<sub>1476</sub> **Interaction between valence and self-reference**

<sub>1477</sub> In this part, we combined the experiments that explicitly manipulated the  
<sub>1478</sub> self-reference and valence, which includes 3a, 3b, 6b, 7a, and 7b. For the positive versus  
<sub>1479</sub> negative contrast, data were from five experiments with 178 participants; for positive  
<sub>1480</sub> versus neutral and neutral versus negative contrasts, data were from three experiments ( <sub>1481</sub> 3a, 3b, and 6b) with 108 participants.

<sub>1482</sub> In most of these experiments, the interaction between self-reference and valence was  
<sub>1483</sub> significant (see results of each experiment in supplementary materials). In the  
<sub>1484</sub> mini-meta-analysis, we analyzed the valence effect for self-referential condition and  
<sub>1485</sub> other-referential condition separately.

1486 For the self-referential condition, we found the same pattern as in the first part of  
1487 results. That is we found significant differences between positive and neutral as well as  
1488 positive and negative, but not neutral and negative. The effect size of RT between positive  
1489 and negative is Cohen's  $d = -0.89 \pm 0.12$ , 95% CI [-1.11 -0.66]; on  $d'$  was Cohen's  $d = 0.61$   
1490  $\pm 0.09$ , 95% CI [0.44 0.78]. The effect was also observed between positive and neutral  
1491 condition, RT: Cohen's  $d = -0.76 \pm 0.13$ , 95% CI [-1.01 -0.50];  $d'$ : Cohen's  $d = 0.69 \pm$   
1492  $0.14$ , 95% CI [0.42 0.96]. And the difference between neutral and bad conditions are not  
1493 significant, RT: Cohen's  $d = 0.03 \pm 0.13$ , 95% CI [-0.22 0.29];  $d'$ : Cohen's  $d = 0.08 \pm 0.08$ ,  
1494 95% CI [-0.07 0.24]. See Figure 18 the middle panel.

1495 For the other-referential condition, we found that only the difference between positive  
1496 and negative on RT was significant, all the other conditions were not. The effect size of RT  
1497 between positive and negative is Cohen's  $d = -0.28 \pm 0.05$ , 95% CI [-0.38 -0.17]; on  $d'$  was  
1498 Cohen's  $d = -0.02 \pm 0.08$ , 95% CI [-0.17 0.13]. The effect was not observed between  
1499 positive and neutral condition, RT: Cohen's  $d = -0.12 \pm 0.10$ , 95% CI [-0.31 0.06];  $d'$ :  
1500 Cohen's  $d = 0.01 \pm 0.08$ , 95% CI [-0.16 0.17]. And the difference between neutral and bad  
1501 conditions are not significant, RT: Cohen's  $d = 0.14 \pm 0.09$ , 95% CI [-0.03 0.31];  $d'$ :  
1502 Cohen's  $d = 0.05 \pm 0.07$ , 95% CI [-0.08 0.18]. See Figure 18 right panel.

### 1503 Generalizability of the valence effect

1504 In this part, we reported the results from experiment 4 in which either moral valence  
1505 or self-reference were manipulated as task-irrelevant stimuli.

1506 For experiment 4a, when self-reference was the target and moral valence was  
1507 task-irrelevant, we found that only under the implicit self-referential condition, i.e., when  
1508 the moral words were presented as task irrelevant stimuli, there was the main effect of  
1509 valence and interaction between valence and reference for both  $d$  prime and RT (See  
1510 supplementary results for the detailed statistics). For  $d$  prime, we found good-self

1511 condition ( $2.55 \pm 0.86$ ) had higher  $d$  prime than bad-self condition ( $2.38 \pm 0.80$ ); good self  
1512 condition was also higher than neutral self ( $2.45 \pm 0.78$ ) but there was not statistically  
1513 significant, while the neutral-self condition was higher than bad self condition and not  
1514 significant neither. For reaction times, good-self condition ( $654.26 \pm 67.09$ ) were faster  
1515 relative to bad-self condition ( $665.64 \pm 64.59$ ), and over neutral-self condition ( $664.26 \pm$   
1516  $64.71$ ). The difference between neutral-self and bad-self conditions were not significant.  
1517 However, for the other-referential condition, there was no significant differences between  
1518 different valence conditions. See Figure 19.

1519 For experiment 4b, when valence was the target and the identity was task-irrelevant,  
1520 we found a strong valence effect (see supplementary results and Figure 20, Figure 21).

1521 In this experiment, the advantage of good-self condition can only be disentangled by  
1522 comparing the self-referential and other-referential conditions. Therefore, we calculated the  
1523 differences between the valence effect under self-referential and other referential conditions  
1524 and used the weighted variance as the variance of this differences. We found this  
1525 modulation effect on RT. The valence effect of RT was stronger in self-referential than  
1526 other-referential for the Good vs. Neutral condition ( $-0.33 \pm 0.01$ ), and to a less extent the  
1527 Good vs. Bad condition ( $-0.17 \pm 0.01$ ). While the size of the other effect's CI included  
1528 zero, suggestion those effects didn't differ from zero. See Figure 22.

### 1529 Specificity of valence effect

1530 In this part, we analyzed the results from experiment 5, which included positive,  
1531 neutral, and negative valence from four different domains: morality, emotion, aesthetics of  
1532 human, and aesthetics of scene. We found interaction between valence and domain for both  
1533  $d$  prime and RT (match trials). A common pattern appeared in all four domains: each  
1534 domain showed a binary results instead of gradient on both  $d$  prime and RT. For morality,  
1535 aesthetics of human, and aesthetics of scene, there was a positivity effect where the positive

1536 conditions had advantages over both neutral (greater  $d$  prime and faster RT), while neutral  
1537 and negative conditions didn't differ from each other. But for the emotional stimuli, there  
1538 was a reversed negativity effect: positive and neutral conditions were not significantly  
1539 different from each other but both had advantage over negative conditions. See  
1540 supplementary materials for detailed statistics. Also note that the effect size in moral  
1541 domain is smaller than the aesthetic domains (beauty of people and beauty of scene). See  
1542 Figure 23.

1543 **Self-reported personal distance**

1544 See Figure 24.

1545 **Correlation analyses**

1546 The reliability of questionnaires can be found in (Liu et al., 2020). We calculated the  
1547 correlation between the data from behavioral task and the questionnaire data. First, we  
1548 calculated the score for each scale based on their structure and factor loading, instead of  
1549 sum score (McNeish & Wolf, 2020). Then, we used SEM to estimate the correlation  
1550 because it can include measurement model and statistical model in a unified framework.

1551 To make sure that what we found were not false positive, we used two method to  
1552 ensure the robustness of our analysis. first, we split the data into two half: the data with  
1553 self and without, then, we used the conditional random forest to find the robust correlation  
1554 in the exploratory data (with self reference) that can be replicated in the confirmatory data  
1555 (without the self reference). The robust correlation were then analyzed using SEM

1556 Instead of use the exploratory correlation analysis, we used a more principled way to  
1557 explore the correlation between parameter of HDM ( $v$ ,  $t$ , and  $a$ ) and scale scores and  
1558 person distance.

1559 We didn't find the correlation between scale scores and the parameters of HDM,

1560 but found weak correlation between personal distance and the parameter estimated from  
1561 Good and neutral conditions.

1562 First, boundary separation (*a*) of moral good condition was correlated with both  
1563 Self-Bad distance ( $r = 0.198$ , 95% CI [],  $p = 0.0063$ ) and Neutral-Bad distance  
1564 ( $r = 0.1571$ , 95% CI [],  $p = 0.031$ ). At the same time, the non-decision time is negatively  
1565 correlated with Self-Bad distance ( $r = 0.169$ , 95% CI [],  $p = 0.0197$ ). See Figure 25.

1566 Second, we found the boundary separation of neutral condition is positively  
1567 correlated with the personal distance between self and good distance ( $r = 0.189$ , 95% CI [],  
1568  $p = 0.036$ ), but negatively correlated with self-neutral distance( $r = -0.183$ , 95% CI [],  
1569  $p = 0.042$ ). Also, the drift rate of the neutral condition is positively correlated with the  
1570 Self-Bad distance ( $r = 0.177$ , 95% CI [],  $p = 0.048$ ).a. See figure 26

1571 We also explored the correlation between behavioral data and questionnaire scores  
1572 separately for experiments with and without self-referential, however, the sample size is  
1573 very low for some conditions.

## 1574 Discussion

## 1575 References

- 1576 Abele, A. E., Ellemers, N., Fiske, S. T., Koch, A., & Yzerbyt, V. (2020). Navigating  
1577 the social world: Toward an integrated framework for evaluating self, individuals,  
1578 and groups. *Psychological Review*. <https://doi.org/10.1037/rev0000262>
- 1579 Amodio, D. M. (2019). Social cognition 2.0: An interactive memory systems  
1580 account. *Trends in Cognitive Sciences*, 23(1), 21–33.  
1581 <https://doi.org/10.1016/j.tics.2018.10.002>
- 1582 Anderson, E., Siegel, E. H., Bliss-Moreau, E., & Barrett, L. F. (2011). The visual  
1583 impact of gossip. *Science*, 332(6036), 1446–1448.

- 1584 https://doi.org/10.1126/science.1201574
- 1585 Brainard, D. H. (1997). The psychophysics toolbox. *Spatial Vision*, 10(4), 433–436.
- 1586 Journal Article.
- 1587 Bürkner, P.-C. (2017). Brms: An r package for bayesian multilevel models using
- 1588 stan. *Journal of Statistical Software; Vol 1, Issue 1 (2017)*. Journal Article.
- 1589 Retrieved from
- 1590 <https://www.jstatsoft.org/v080/i01%20http://dx.doi.org/10.18637/jss.v080.i01>
- 1591 Carlson, R. W., Maréchal, M. A., Oud, B., Fehr, E., & Crockett, M. J. (2020).
- 1592 Motivated misremembering of selfish decisions. *Nature Communications*, 11(1),
- 1593 2100. <https://doi.org/10.1038/s41467-020-15602-4>
- 1594 Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M.,
- 1595 ... Riddell, A. (2017). Stan: A probabilistic programming language. *Journal of*
- 1596 *Statistical Software*, 76(1). Journal Article.
- 1597 <https://doi.org/10.18637/jss.v076.i01>
- 1598 Dang, J., King, K. M., & Inzlicht, M. (2020). Why are self-report and behavioral
- 1599 measures weakly correlated? *Trends in Cognitive Sciences*, 24(4), 267–269.
- 1600 <https://doi.org/10.1016/j.tics.2020.01.007>
- 1601 Dzhelyova, M., Perrett, D. I., & Jentzsch, I. (2012). Temporal dynamics of
- 1602 trustworthiness perception. *Brain Research*, 1435, 81–90.
- 1603 <https://doi.org/10.1016/j.brainres.2011.11.043>
- 1604 Ellemers, N., Toorn, J. van der, Paunov, Y., & Leeuwen, T. van. (2019). The
- 1605 psychology of morality: A review and analysis of empirical studies published
- 1606 from 1940 through 2017. *Personality and Social Psychology Review*, 23(4),
- 1607 332–366. <https://doi.org/10.1177/1088868318811759>
- 1608 Enock, F. E., Hewstone, M. R. C., Lockwood, P. L., & Sui, J. (2020). Overlap in

- 1609 processing advantages for minimal ingroups and the self. *Scientific Reports*,  
1610 10(1), 18933. <https://doi.org/10.1038/s41598-020-76001-9>
- 1611 Enock, F., Sui, J., Hewstone, M., & Humphreys, G. W. (2018). Self and team  
1612 prioritisation effects in perceptual matching: Evidence for a shared  
1613 representation. *Acta Psychologica*, 182, 107–118.  
1614 <https://doi.org/10.1016/j.actpsy.2017.11.011>
- 1615 Faul, F., Erdfelder, E., Buchner, A., & Lang, A.-G. (2009). Statistical power  
1616 analyses using g\*power 3.1: Tests for correlation and regression analyses.  
1617 *Behavior Research Methods*, 41(4), 1149–1160.  
1618 <https://doi.org/10.3758/BRM.41.4.1149>
- 1619 Firestone, C., & Scholl, B. J. (2015). Enhanced visual awareness for morality and  
1620 pajamas? Perception vs. Memory in ‘top-down’ effects. *Cognition*, 136, 409–416.  
1621 <https://doi.org/10.1016/j.cognition.2014.10.014>
- 1622 Freeman, J. B., & Ambady, N. (2011). A dynamic interactive theory of person  
1623 construal. *Psychological Review*, 118(2), 247–279.  
1624 <https://doi.org/10.1037/a0022327>
- 1625 Gantman, A. P., & Van Bavel, J. J. (2014). The moral pop-out effect: Enhanced  
1626 perceptual awareness of morally relevant stimuli. *Cognition*, 132(1), 22–29.  
1627 <https://doi.org/10.1016/j.cognition.2014.02.007>
- 1628 Goh, J. X., Hall, J. A., & Rosenthal, R. (2016). Mini meta-analysis of your own  
1629 studies: Some arguments on why and a primer on how. *Social and Personality  
1630 Psychology Compass*, 10(10), 535–549. Journal Article.  
1631 <https://doi.org/10.1111/spc3.12267>
- 1632 Goodwin, G. P. (2015). Moral character in person perception. *Current Directions in  
1633 Psychological Science*, 24(1), 38–44. <https://doi.org/10.1177/0963721414550709>

- 1634 Goodwin, G. P., Piazza, J., & Rozin, P. (2014). Moral character predominates in  
1635 person perception and evaluation. *Journal of Personality and Social Psychology*,  
1636 106(1), 148–168. <https://doi.org/10.1037/a0034726>
- 1637 Henrich, J., Heine, S. J., & Norenzayan, A. (2010). The weirdest people in the  
1638 world? *Behavioral and Brain Sciences*, 33(2), 61–83.  
1639 <https://doi.org/10.1017/S0140525X0999152X>
- 1640 Hofmann, W., Wisneski, D. C., Brandt, M. J., & Skitka, L. J. (2014). Morality in  
1641 everyday life. *Science*, 345(6202), 1340–1343.  
1642 <https://doi.org/10.1126/science.1251560>
- 1643 Hu, C.-P., Lan, Y., Macrae, C. N., & Sui, J. (2020). Good me bad me: Does valence  
1644 influence self-prioritization during perceptual decision-making? *Collabra:*  
1645 *Psychology*, 6(1), 20. Journal Article. <https://doi.org/10.1525/collabra.301>
- 1646 Jones, B. C., DeBruine, L. M., Flake, J. K., Liuzza, M. T., Antfolk, J., Arinze, N.  
1647 C., ... Coles, N. A. (2020). To which world regions does the valence-dominance  
1648 model of social perception apply? *Nature Human Behaviour*.  
1649 <https://doi.org/10.31234/osf.io/n26dy>
- 1650 Juechems, K., & Summerfield, C. (2019). Where does value come from? *Trends in*  
1651 *Cognitive Sciences*, 23(10), 836–850. <https://doi.org/10.1016/j.tics.2019.07.012>
- 1652 Kirby, K. N., & Gerlanc, D. (2013). BootES: An r package for bootstrap confidence  
1653 intervals on effect sizes. *Behavior Research Methods*, 45(4), 905–927.  
1654 <https://doi.org/10.3758/s13428-013-0330-5>
- 1655 Klein, N., & Epley, N. (2016). Maybe holier, but definitely less evil, than you:  
1656 Bounded self-righteousness in social judgment. *Journal of Personality and Social*  
1657 *Psychology*, 110(5), 660–674. <https://doi.org/10.1037/pspa0000050>
- 1658 Liu, Q., Wang, F., Yan, W., Peng, K., Sui, J., & Hu, C.-P. (2020). Questionnaire

- 1659 data from the revision of a chinese version of free will and determinism plus  
1660 scale. *Journal of Open Psychology Data*, 8(1), 1. Journal Article.  
1661 <https://doi.org/10.5334/jopd.49/>
- 1662 Matzke, D., & Wagenmakers, E.-J. (2009). Psychological interpretation of the  
1663 ex-gaussian and shifted wald parameters: A diffusion model analysis.  
1664 *Psychonomic Bulletin & Review*, 16(5), 798–817.  
1665 <https://doi.org/10.3758/PBR.16.5.798>
- 1666 McHugh, C., McGann, M., Igou, E. R., & Kinsella, E. (2019). *Moral judgment as*  
1667 *categorization (MJAC)*. PsyArXiv. <https://doi.org/10.31234/osf.io/72dzp>
- 1668 McNeish, D., & Wolf, M. G. (2020). Thinking twice about sum scores. *Behavior*  
1669 *Research Methods*. <https://doi.org/10.3758/s13428-020-01398-0>
- 1670 Monin, B., & Jordan, A. H. (2009). The dynamic moral self: A social psychological  
1671 perspective. In *Personality, identity, and character: Explorations in moral*  
1672 *psychology* (pp. 341–354). New York, NY, US: Cambridge University Press.  
1673 <https://doi.org/10.1017/CBO9780511627125.016>
- 1674 Pelli, D. G. (1997). The VideoToolbox software for visual psychophysics:  
1675 Transforming numbers into movies. *Spatial Vision*, 10(4), 437–442. Journal  
1676 Article.
- 1677 Reicher, S., & Hopkins, N. (2016). Perception, action, and the social dynamics of  
1678 the variable self. *Psychological Inquiry*, 27(4), 341–347.  
1679 <https://doi.org/10.1080/1047840X.2016.1217584>
- 1680 Rouder, J. N., & Lu, J. (2005). An introduction to bayesian hierarchical models  
1681 with an application in the theory of signal detection. *Psychonomic Bulletin &*  
1682 *Review*, 12(4), 573–604. Journal Article. <https://doi.org/10.3758/bf03196750>
- 1683 Rousselet, G. A., & Wilcox, R. R. (2019). Reaction times and other skewed

- 1684 distributions: Problems with the mean and the median. *Meta-Psychology*.  
1685 preprint. <https://doi.org/10.1101/383935>
- 1686 Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2013). Life after p-hacking.  
1687 Conference Proceedings. <https://doi.org/10.2139/ssrn.2205186>
- 1688 Stanley, M. L., & De Brigard, F. (2019). Moral memories and the belief in the good  
1689 self. *Current Directions in Psychological Science*, 28(4), 387–391.  
1690 <https://doi.org/10.1177/0963721419847990>
- 1691 Stein, T., Grubb, C., Bertrand, M., Suh, S. M., & Verosky, S. C. (2017). No impact  
1692 of affective person knowledge on visual awareness: Evidence from binocular  
1693 rivalry and continuous flash suppression. *Emotion*, 17(8), 1199–1207.  
1694 <https://doi.org/10.1037/emo0000305>
- 1695 Stolier, R. M., & Freeman, J. B. (2016). Functional and temporal considerations for  
1696 top-down influences in social perception. *Psychological Inquiry*, 27(4), 352–357.  
1697 <https://doi.org/10.1080/1047840X.2016.1216034>
- 1698 Strohminger, N., Knobe, J., & Newman, G. (2017). The true self: A psychological  
1699 concept distinct from the self: *Perspectives on Psychological Science*.  
1700 <https://doi.org/10.1177/1745691616689495>
- 1701 Sui, J., He, X., & Humphreys, G. W. (2012). Perceptual effects of social salience:  
1702 Evidence from self-prioritization effects on perceptual matching. *Journal of  
1703 Experimental Psychology: Human Perception and Performance*, 38(5),  
1704 1105–1117. Journal Article. <https://doi.org/10.1037/a0029792>
- 1705 Tappin, B. M., & McKay, R. T. (2017). The illusion of moral superiority. *Social  
1706 Psychological and Personality Science*, 8(6), 623–631.  
1707 <https://doi.org/10.1177/1948550616673878>
- 1708 Turner, J. C., Hogg, M. A., Oakes, P. J., Reicher, S. D., & Wetherell, M. S. (1987).

- 1709           *Rediscovering the social group: A self-categorization theory.* Cambridge, MA,  
1710           US: Basil Blackwell.
- 1711           Turner, J. C., Oakes, P. J., Haslam, S. A., & McGarty, C. (1994). Self and  
1712           collective: Cognition and social context. *Personality and Social Psychology  
1713           Bulletin*, 20(5), 454–463. <https://doi.org/10.1177/0146167294205002>
- 1714           Uhlmann, E. L., Pizarro, D. A., & Diermeier, D. (2015). A person-centered  
1715           approach to moral judgment: *Perspectives on Psychological Science*.  
1716           <https://doi.org/10.1177/1745691614556679>
- 1717           Verosky, S. C., & Todorov, A. (2010). Differential neural responses to faces  
1718           physically similar to the self as a function of their valence. *NeuroImage*, 49(2),  
1719           1690–1698. <https://doi.org/10.1016/j.neuroimage.2009.10.017>
- 1720           Waytz, A., Dungan, J., & Young, L. (2013). The whistleblower's dilemma and the  
1721           fairness–loyalty tradeoff. *Journal of Experimental Social Psychology*, 49(6),  
1722           1027–1033. <https://doi.org/10.1016/j.jesp.2013.07.002>
- 1723           Wiecki, T. V., Sofer, I., & Frank, M. J. (2013). HDDM: Hierarchical bayesian  
1724           estimation of the drift-diffusion model in python. *Frontiers in Neuroinformatics*,  
1725           7. <https://doi.org/10.3389/fninf.2013.00014>
- 1726           Willis, J., & Todorov, A. (2006). First impressions: Making up your mind after a  
1727           100-ms exposure to a face. *Psychological Science*, 17(7), 592–598.  
1728           <https://doi.org/10.1111/j.1467-9280.2006.01750.x>
- 1729           Xiao, Y. J., Coppin, G., & Bavel, J. J. V. (2016). Perceiving the world through  
1730           group-colored glasses: A perceptual model of intergroup relations. *Psychological  
1731           Inquiry*, 27(4), 255–274. <https://doi.org/10.1080/1047840X.2016.1199221>

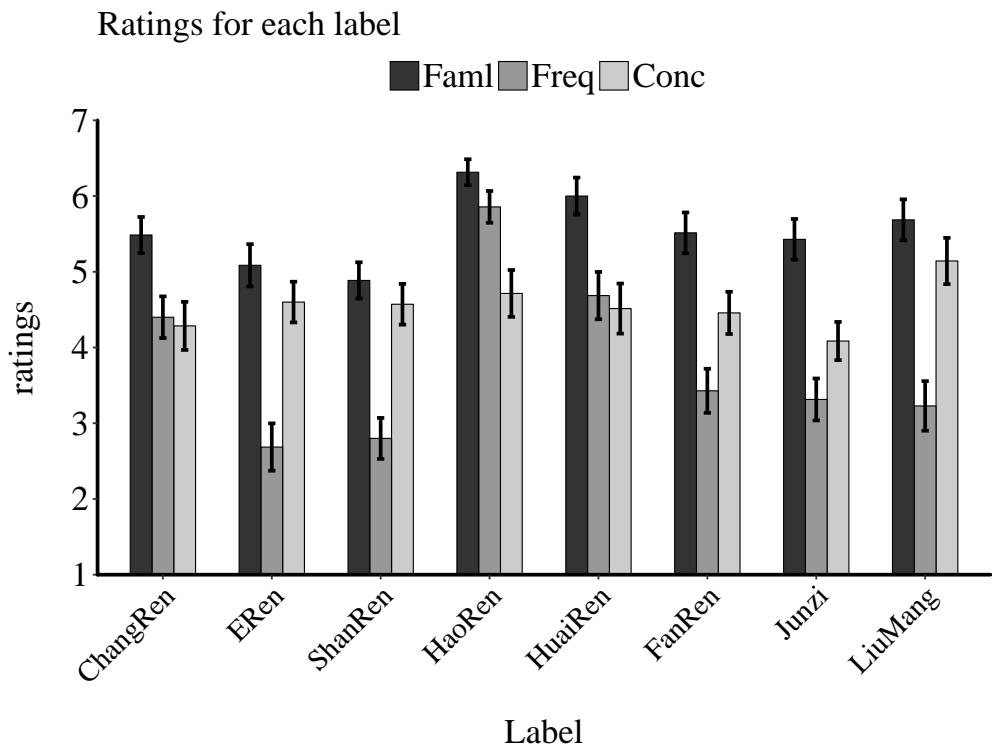


Figure 1. Ratings of words in exp 1b

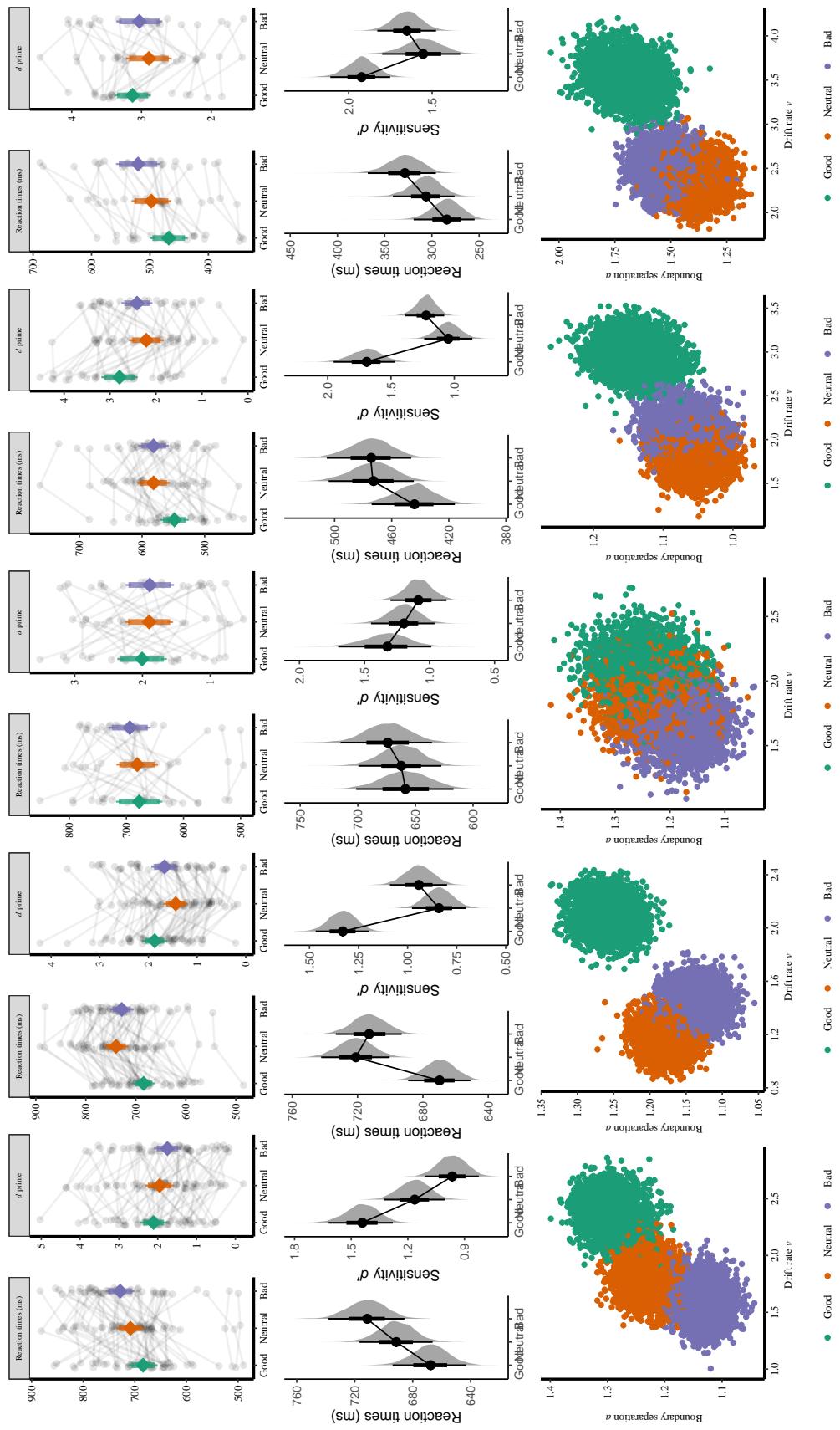


Figure 2. Results for part 1.

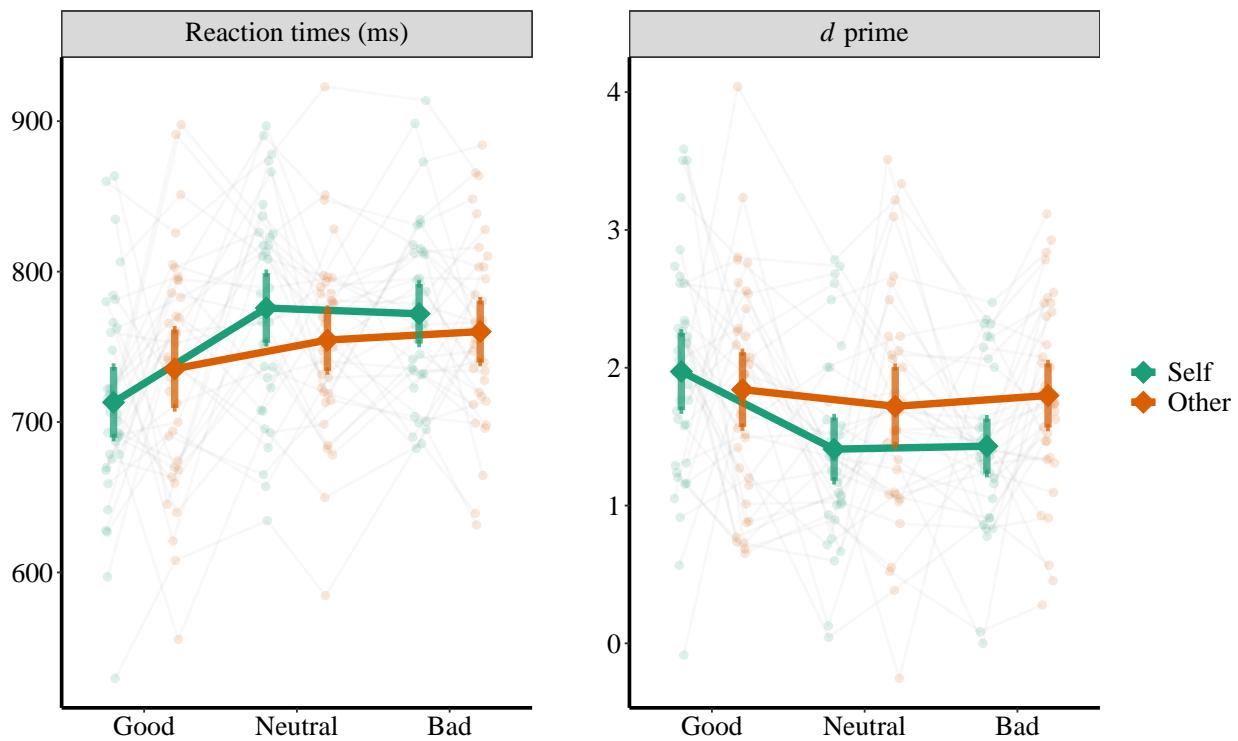


Figure 3. RT and  $d'$  of Experiment 3a.

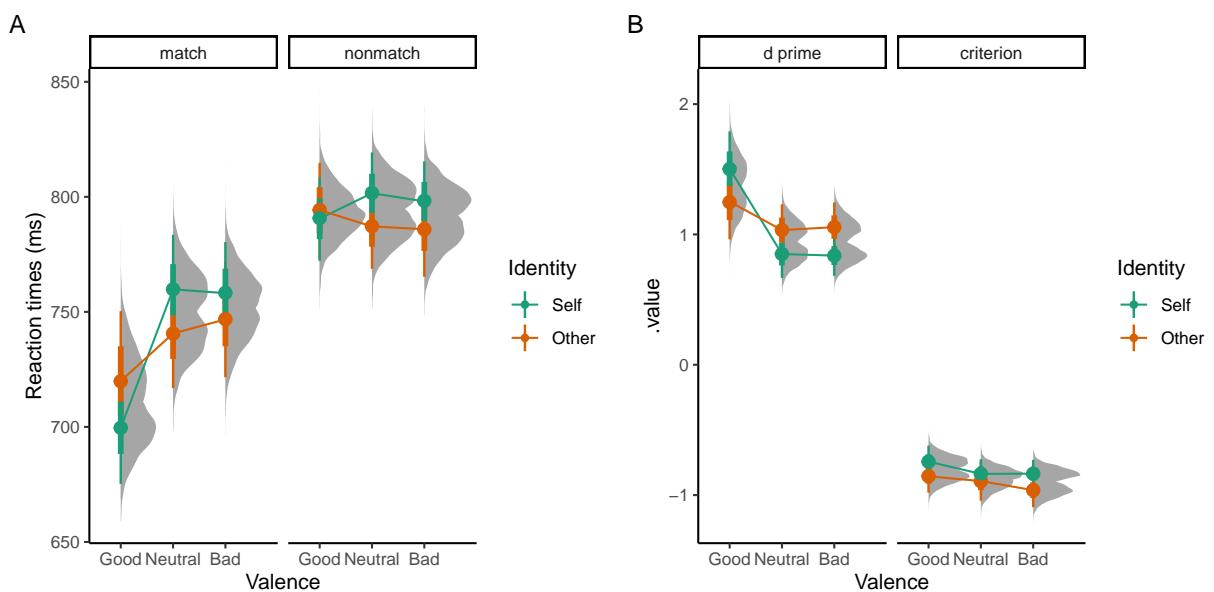


Figure 4. Exp3a: Results of Bayesian GLM analysis.

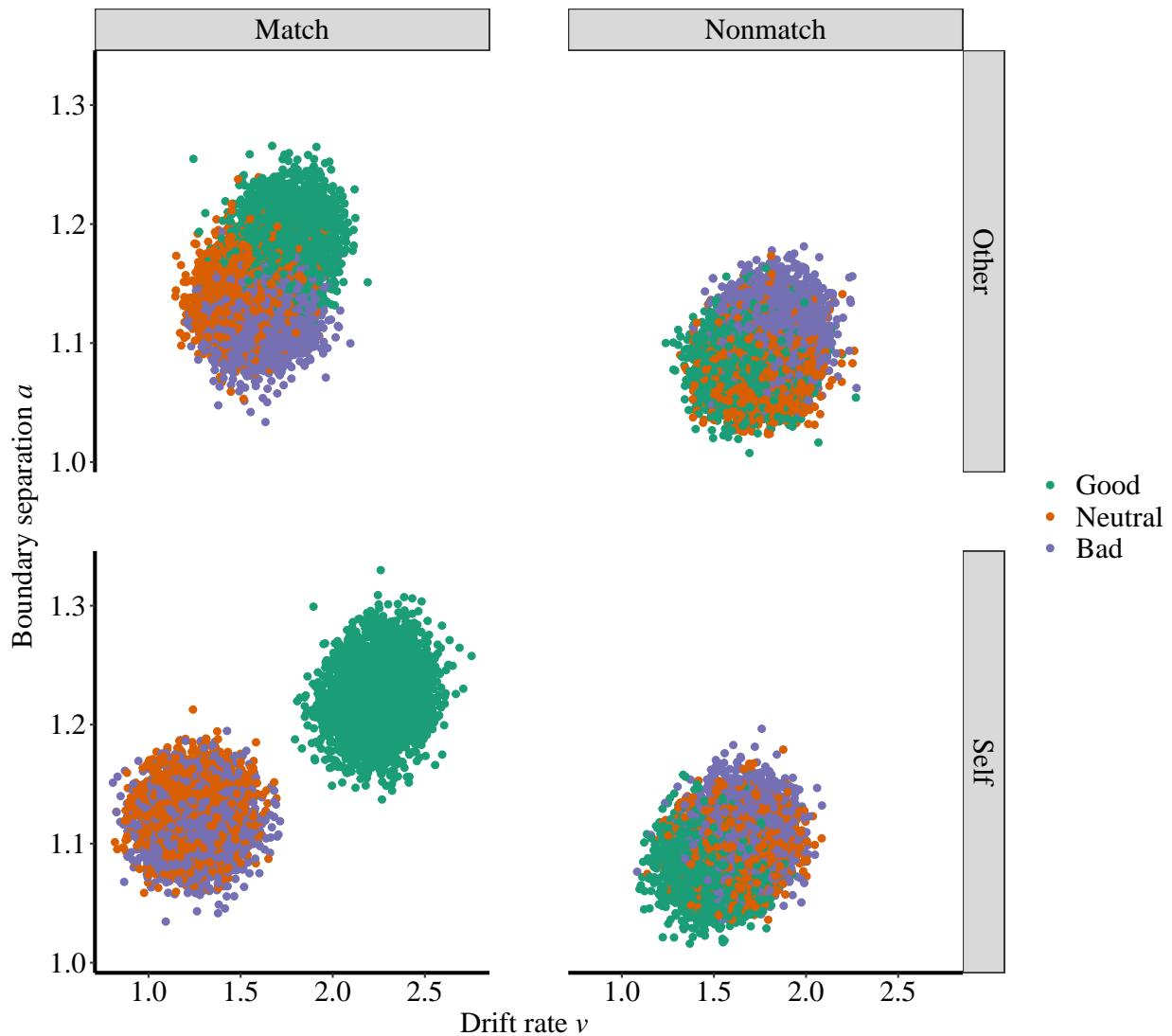


Figure 5. Exp3a: Results of HDDM.

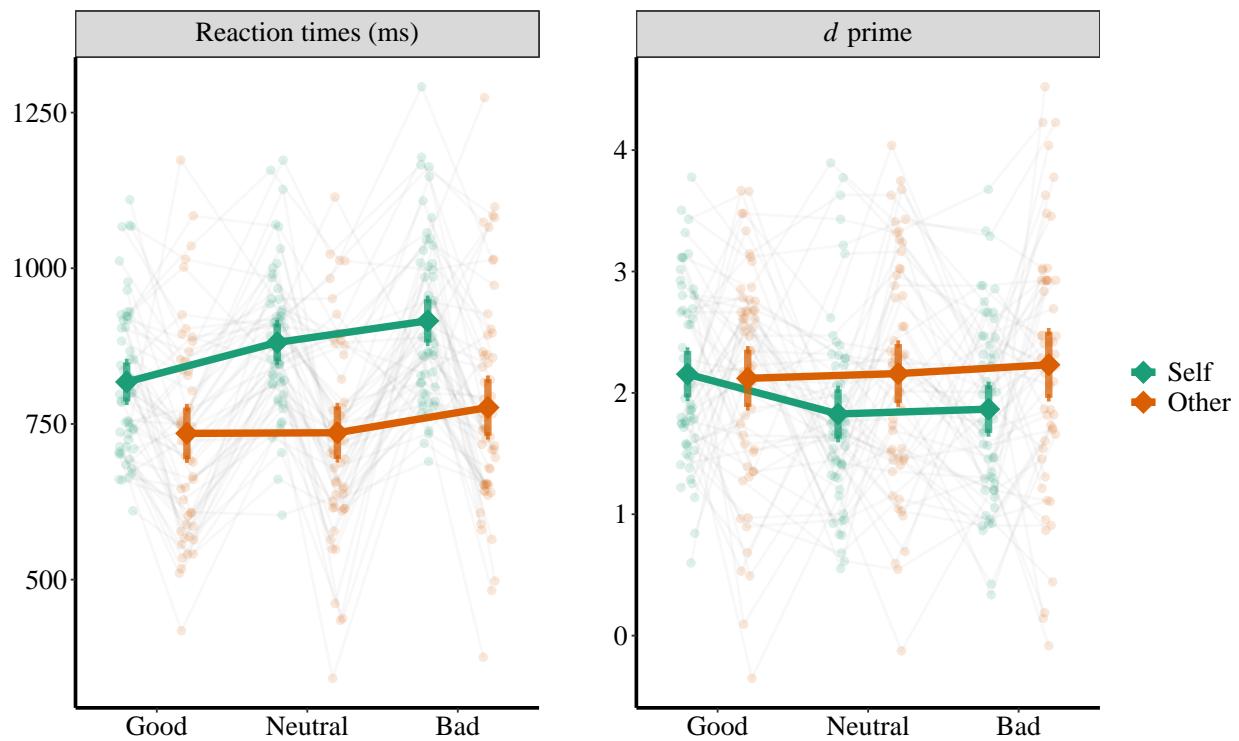
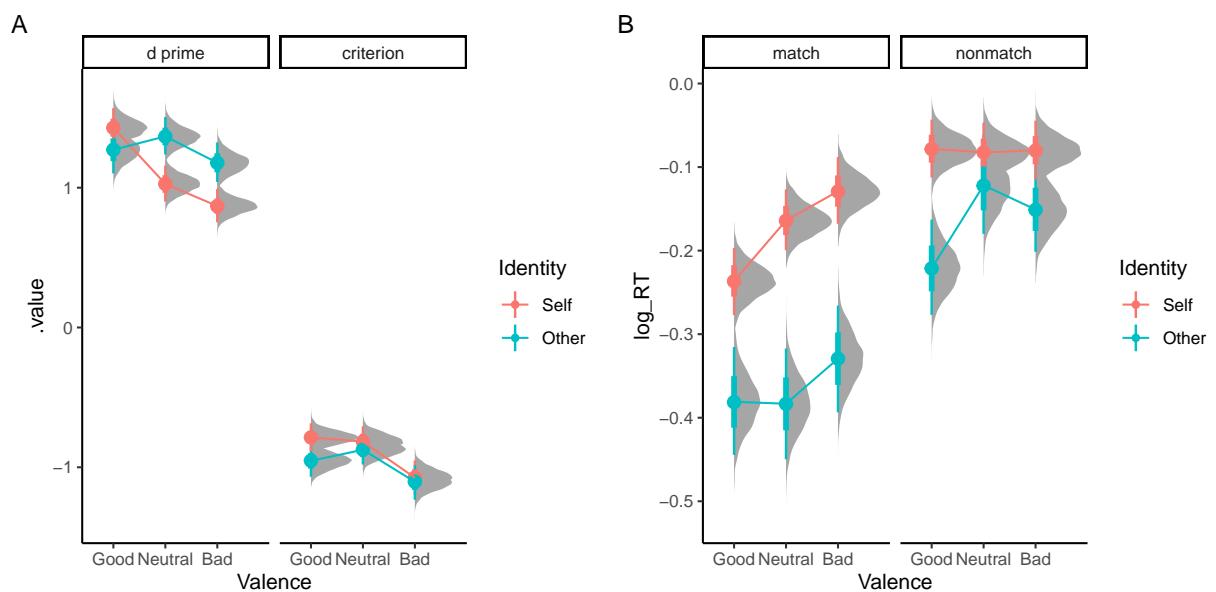
Figure 6. RT and  $d'$  of Experiment 3b.

Figure 7. exp3b: Results of Bayesian GLM analysis.

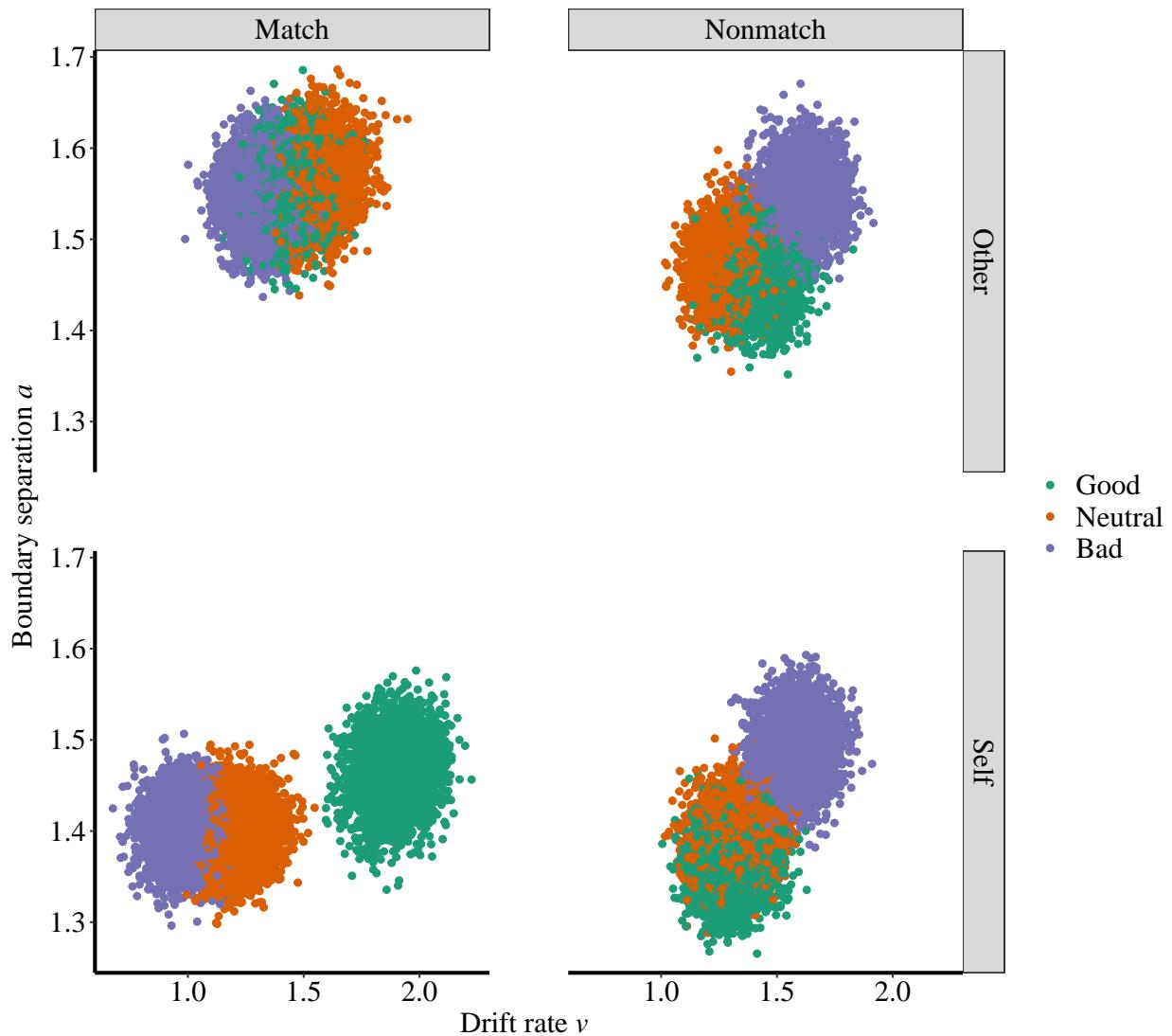


Figure 8. exp3b: Results of HDDM.

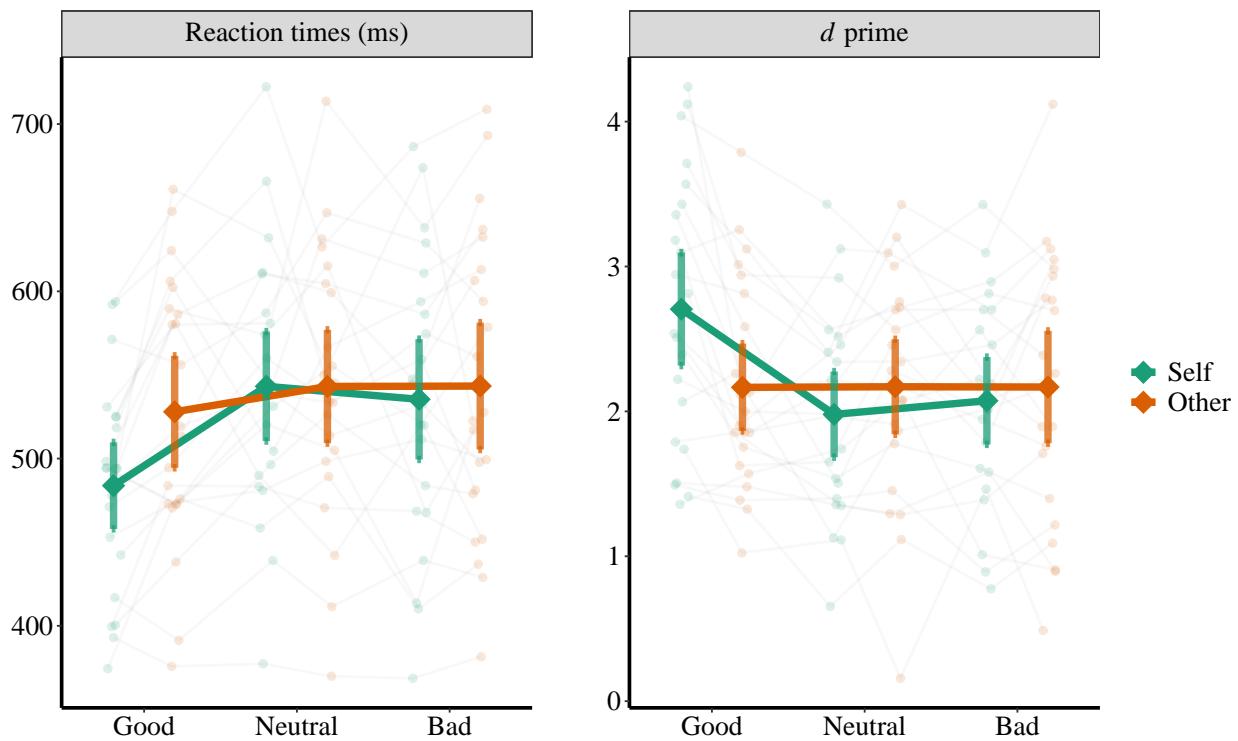


Figure 9. RT and  $d'$  of Experiment 6b.

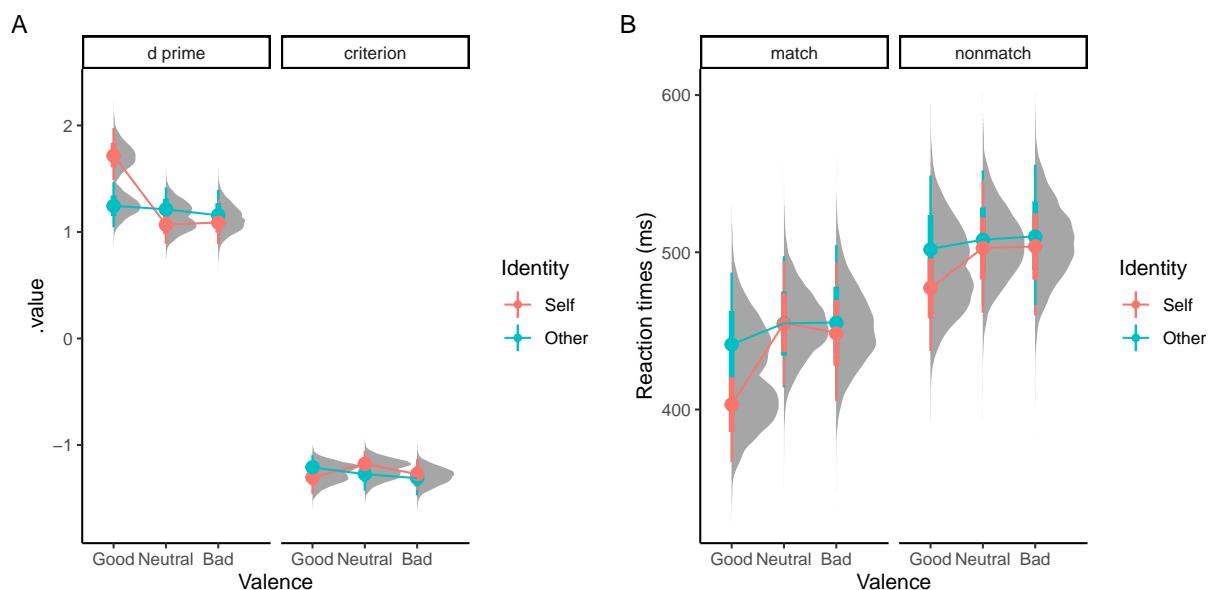


Figure 10. exp6b\_d1: Results of Bayesian GLM analysis.

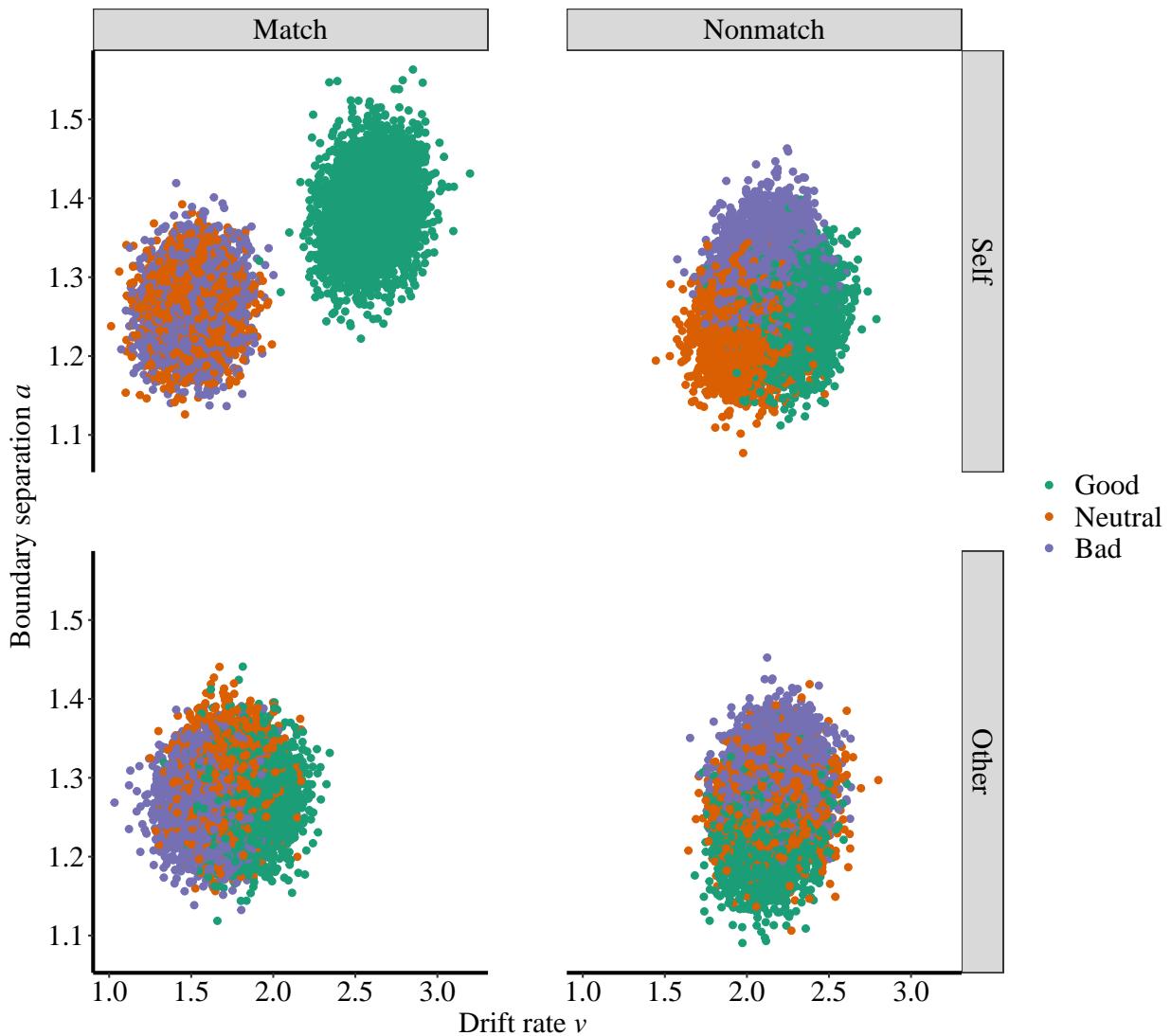


Figure 11. exp6b: Results of HDDM (Day 1).

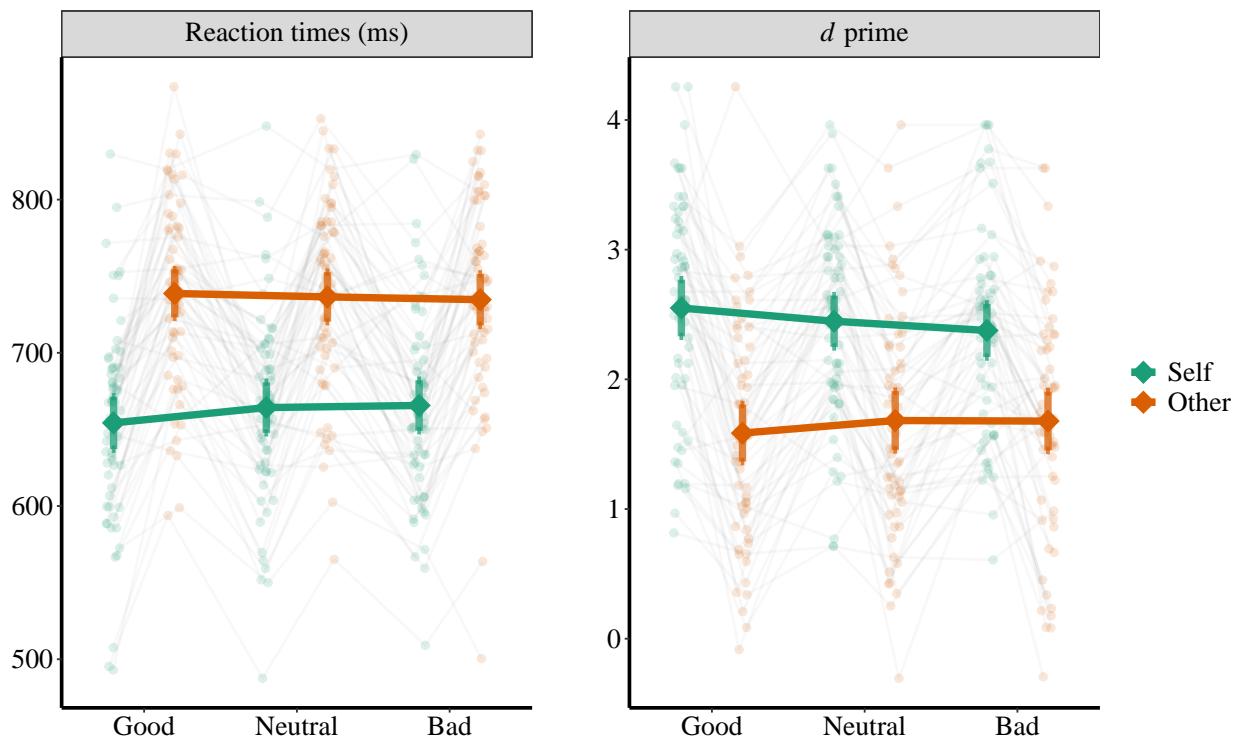


Figure 12. RT and  $d'$  of Experiment 4a.

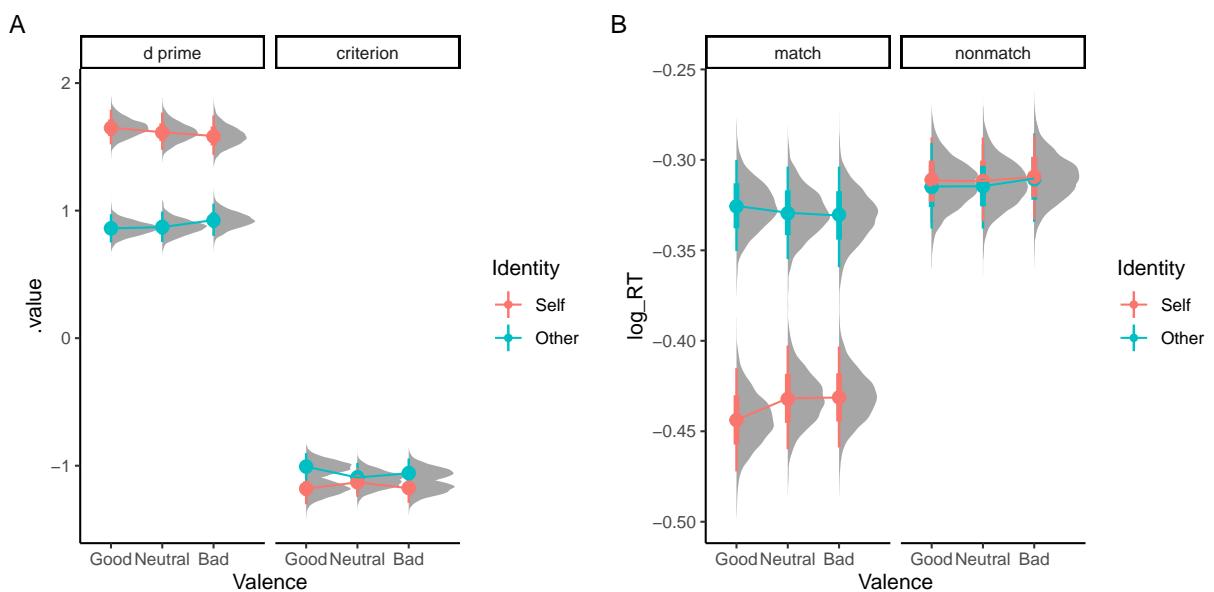


Figure 13. exp4a: Results of Bayesian GLM analysis.

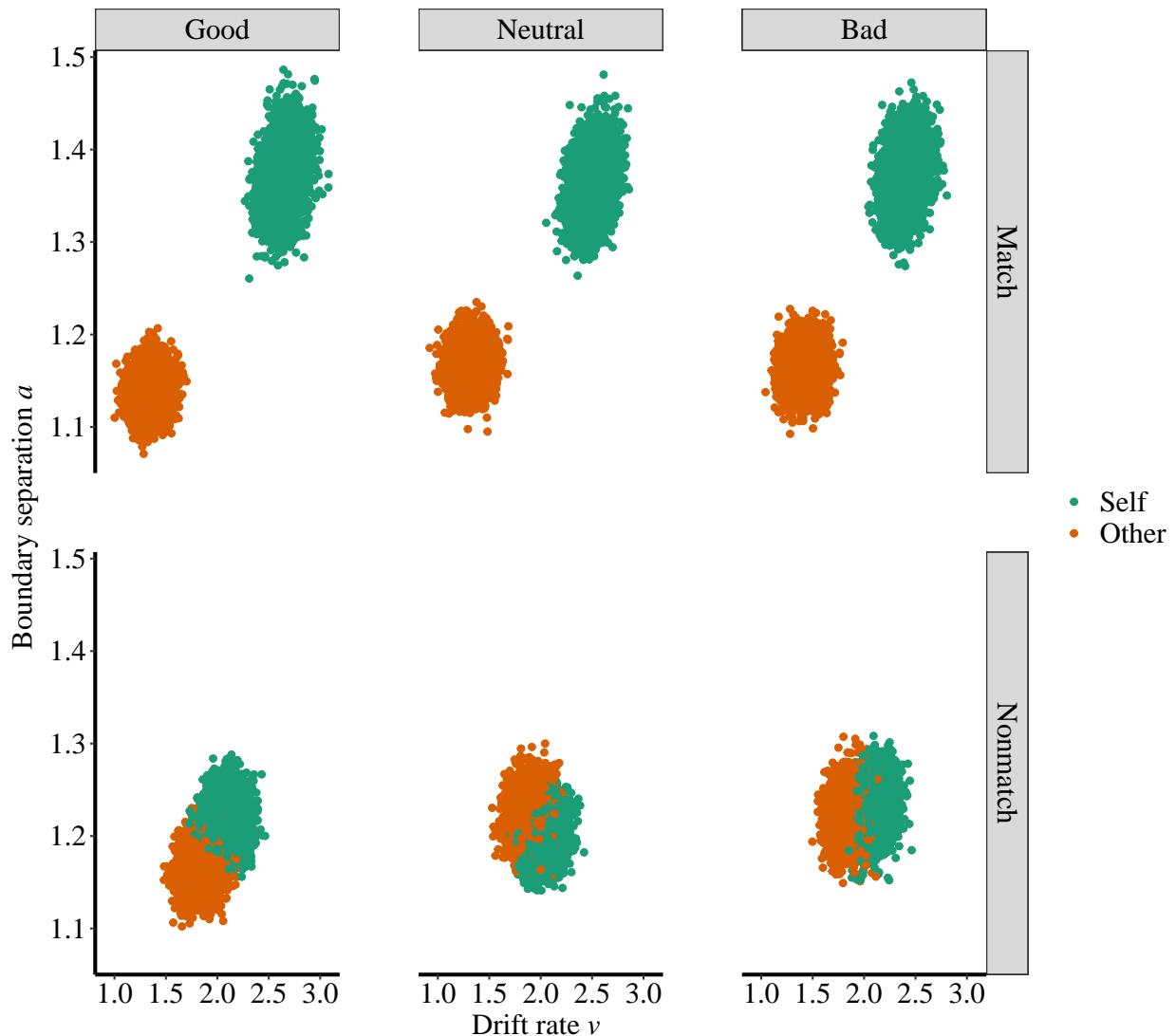


Figure 14. exp4a: Results of HDDM.

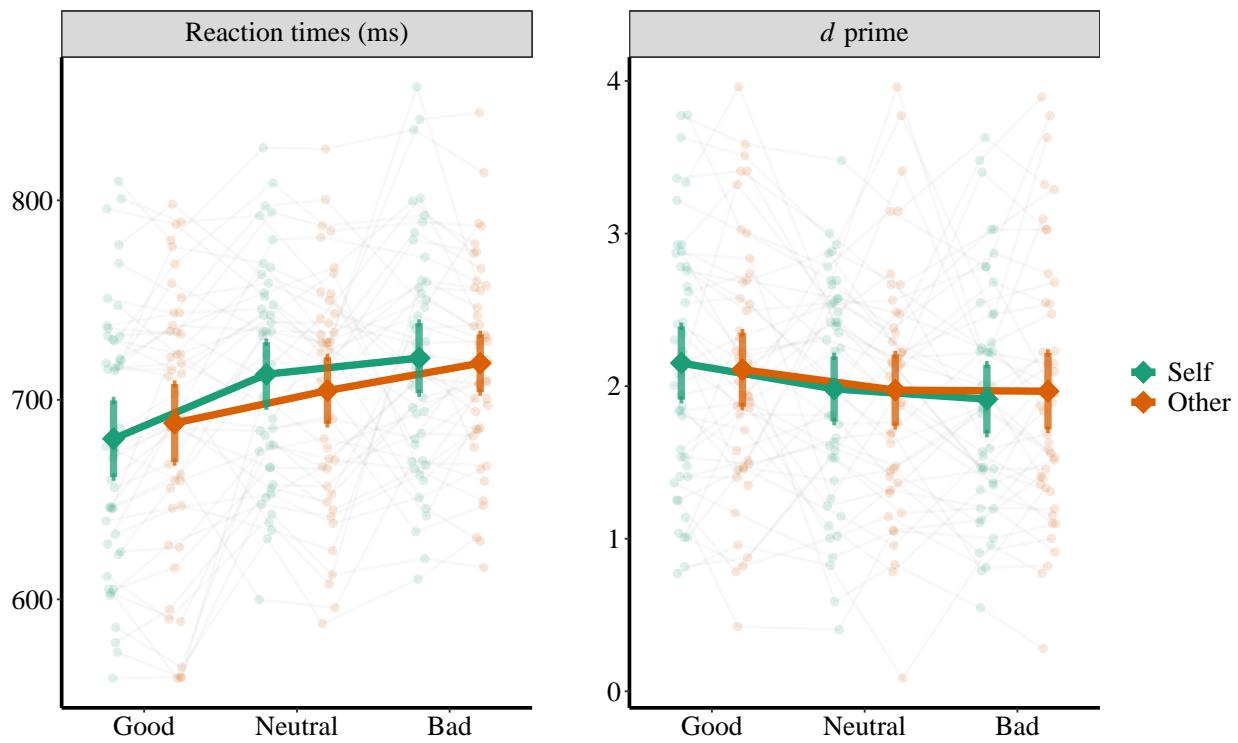


Figure 15. RT and  $d'$  prime of Experiment 4b.

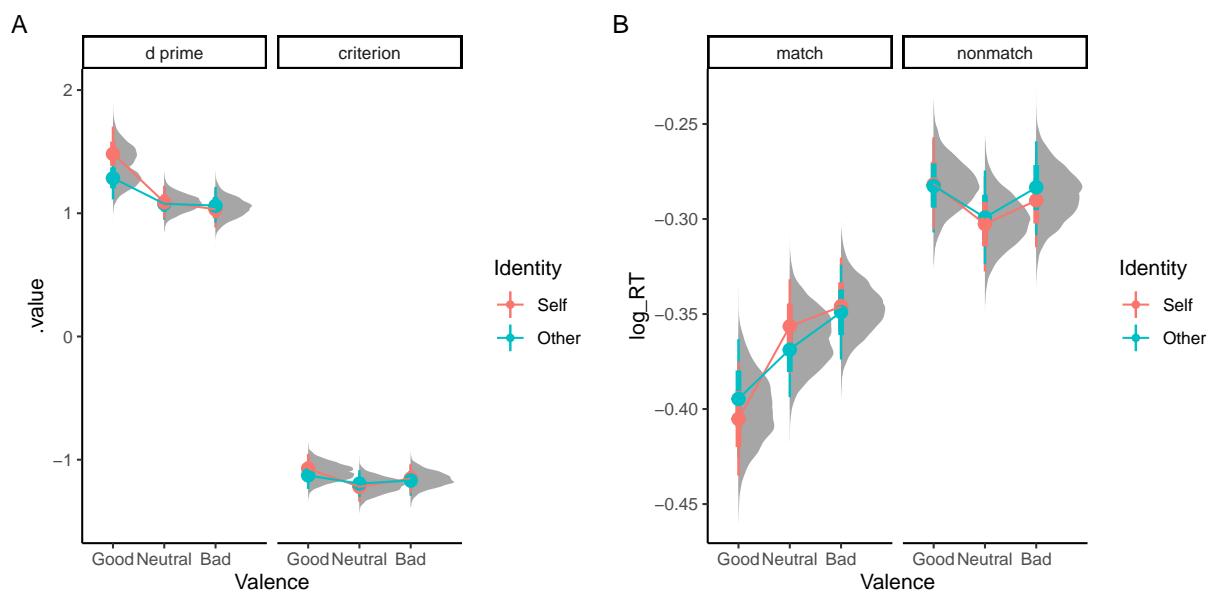


Figure 16. exp4b: Results of Bayesian GLM analysis.

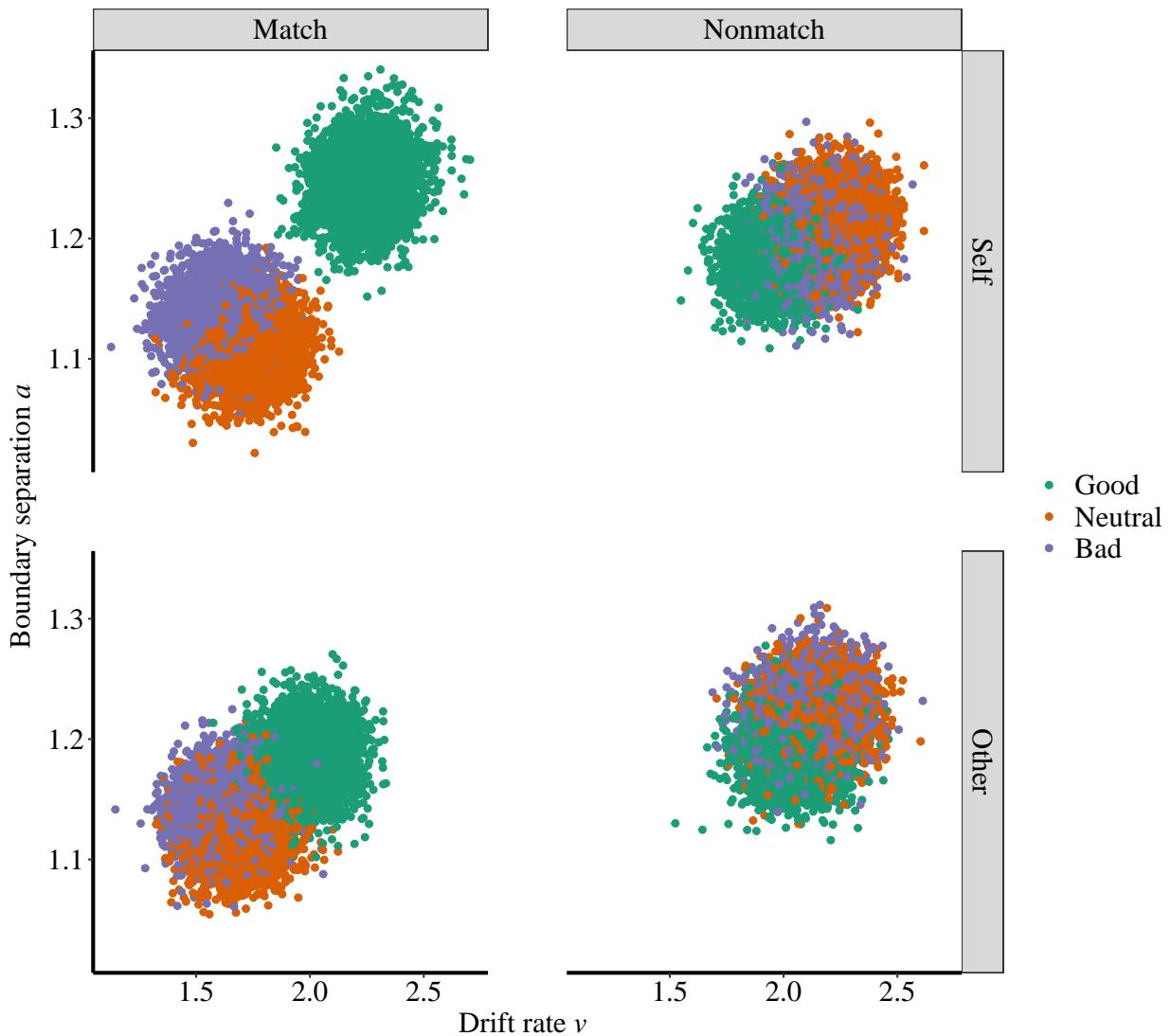


Figure 17. exp4b: Results of HDDM.

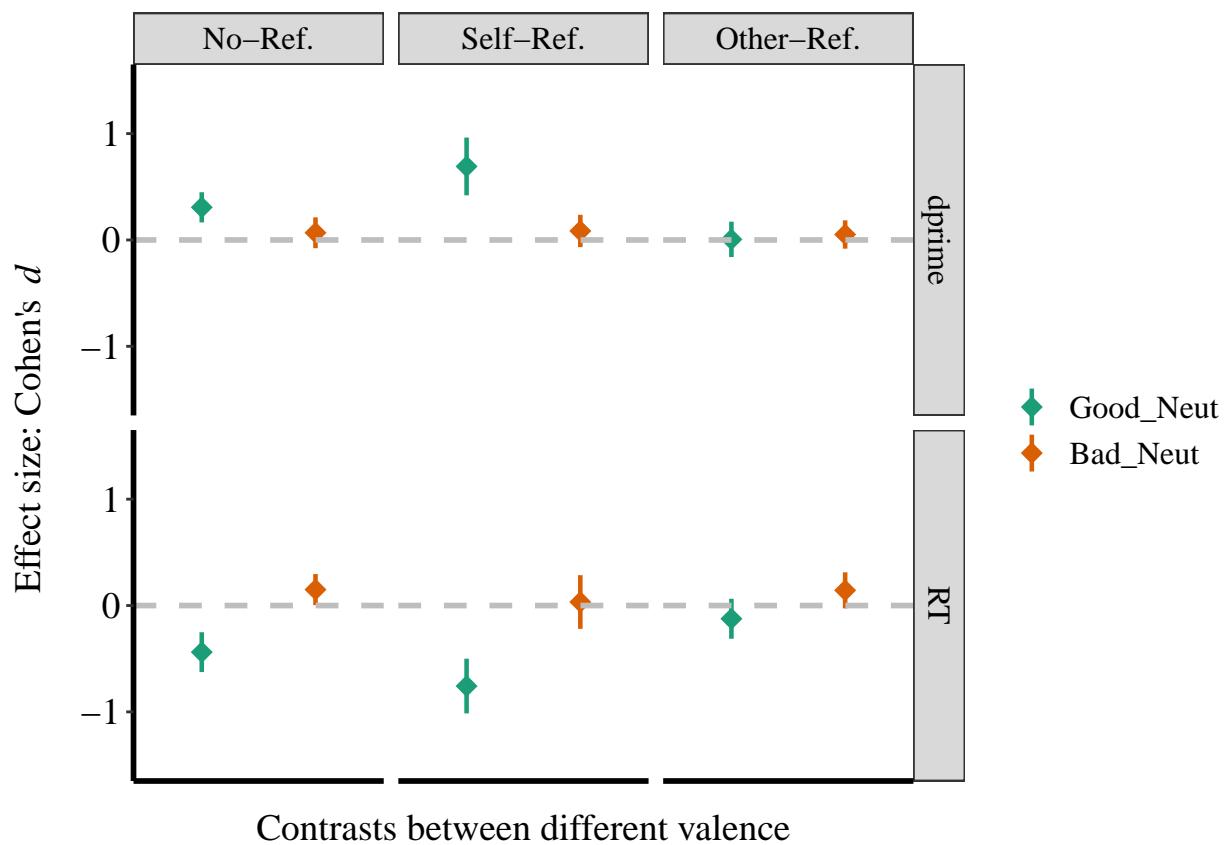


Figure 18. Effect size (Cohen's  $d$ ) of Valence.

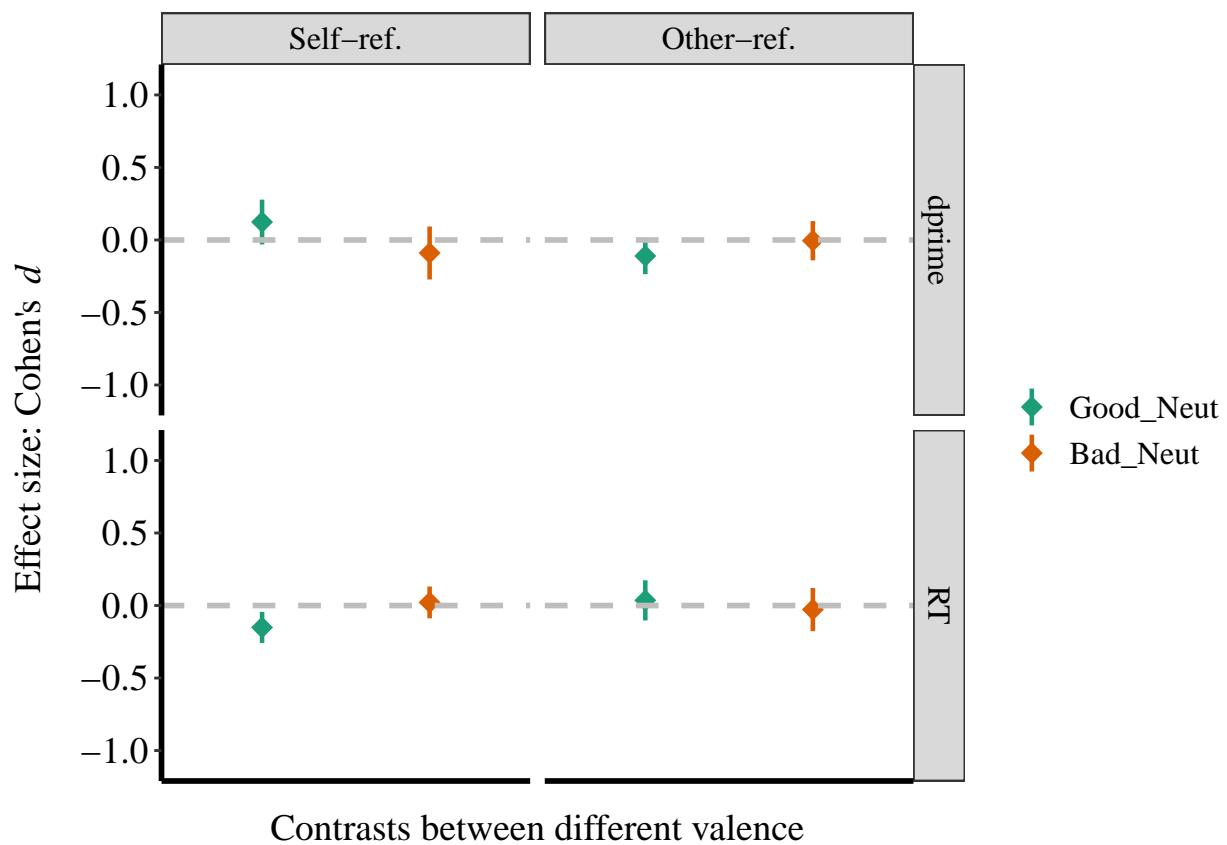


Figure 19. Effect size (Cohen's  $d$ ) of Valence in Exp4a.

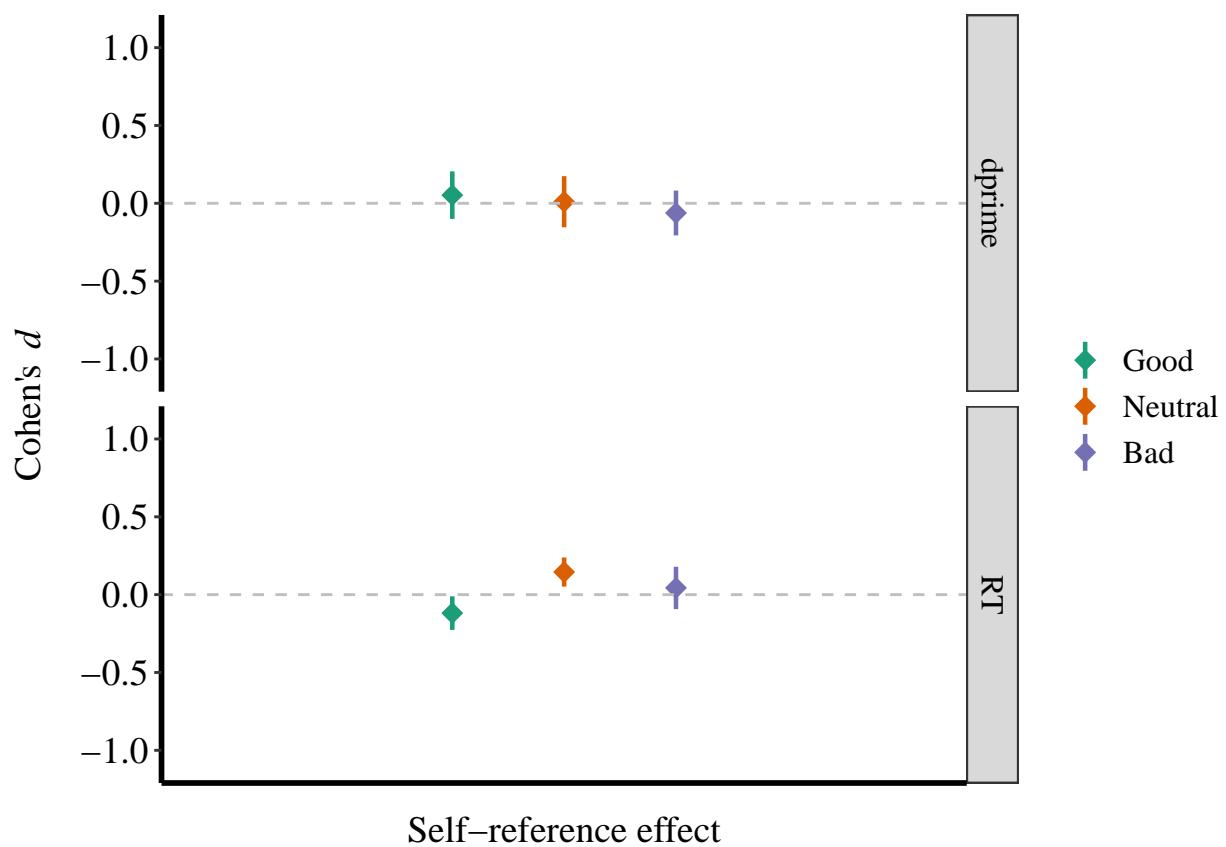


Figure 20. Effect size (Cohen's  $d$ ) of Valence in Exp4b.

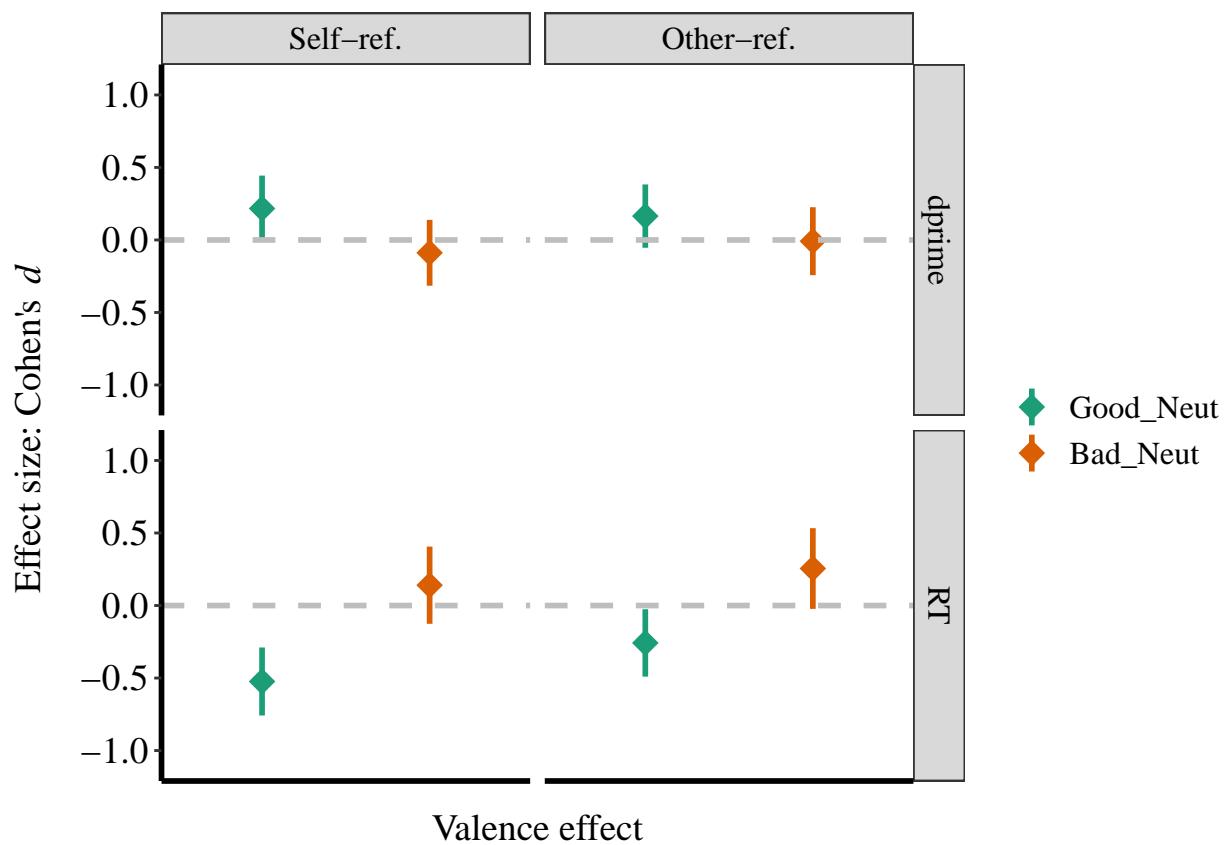


Figure 21. Effect size (Cohen's  $d$ ) of Valence in Exp4b.

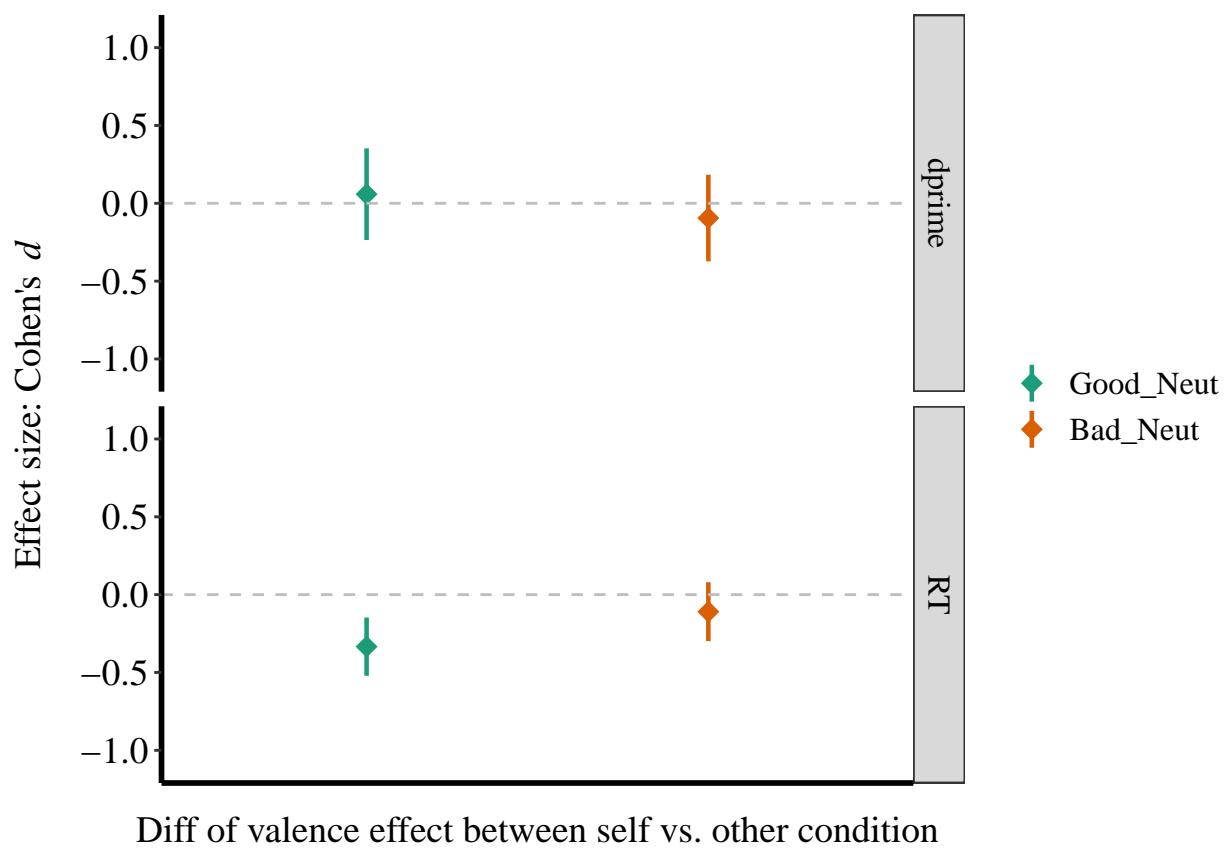


Figure 22. Effect size (Cohen's  $d$ ) of Valence in Exp4b.

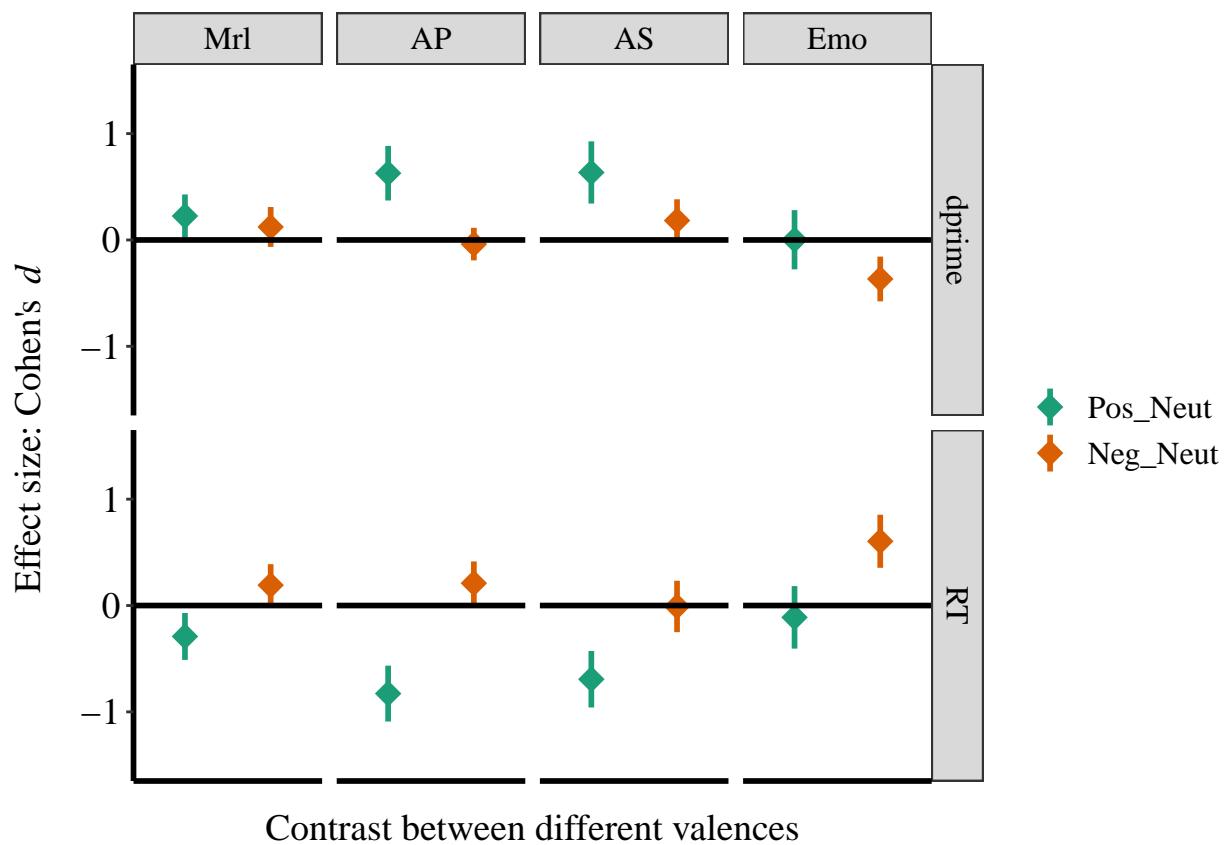


Figure 23. Effect size (Cohen's  $d$ ) of Valence in Exp5.

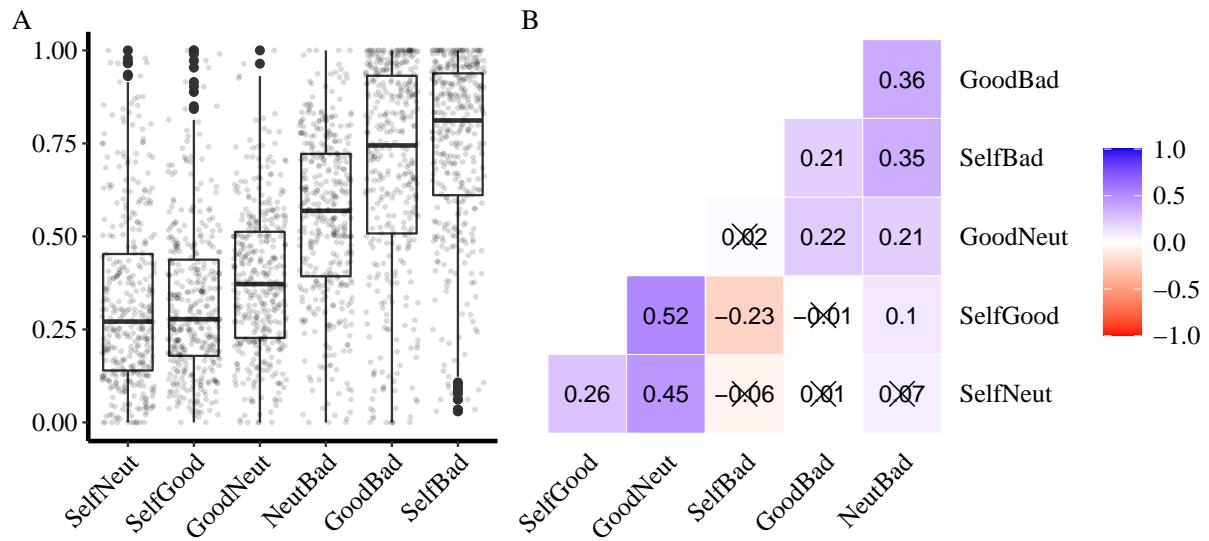


Figure 24. Self-rated personal distance

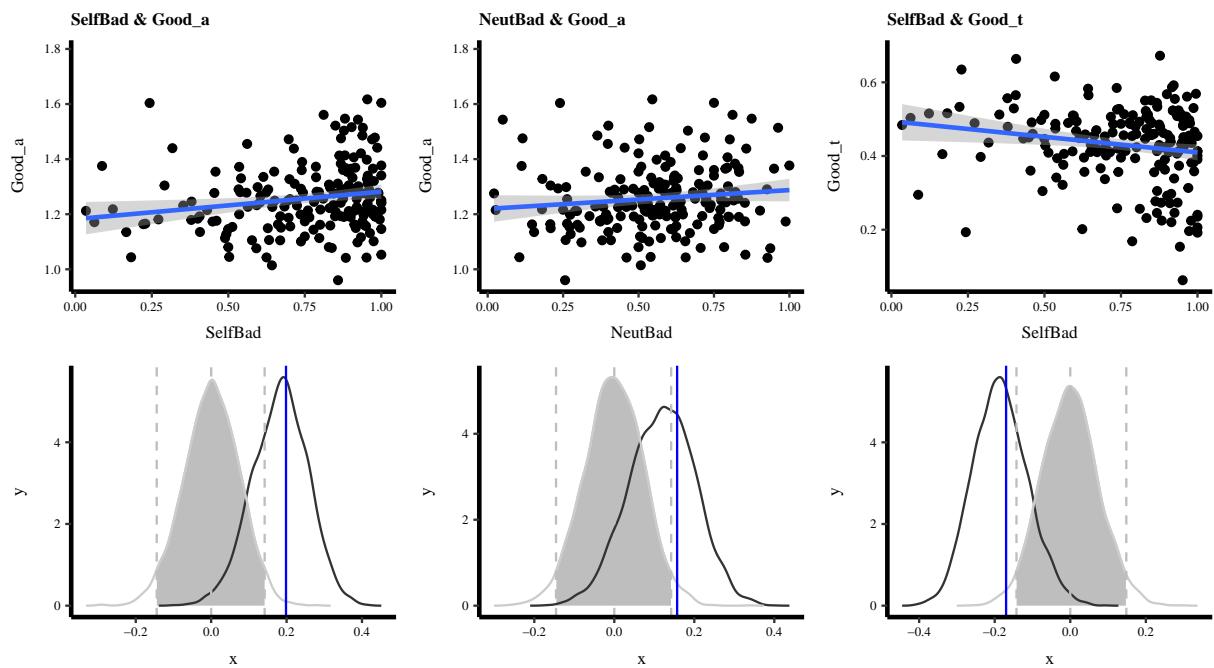
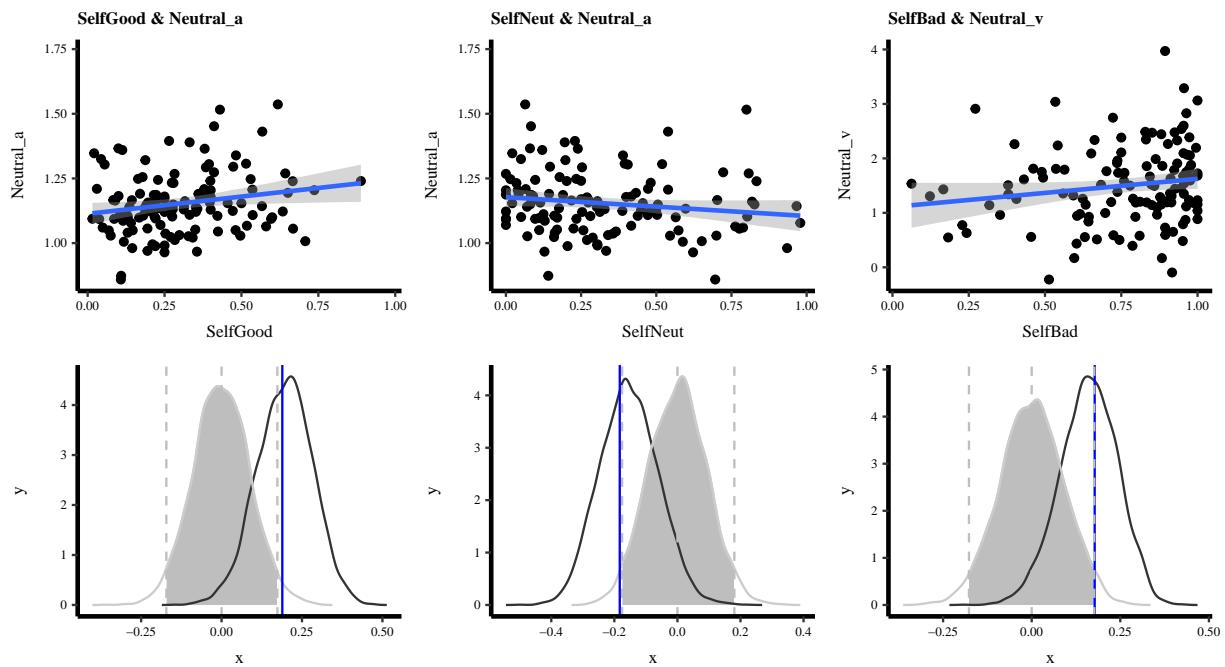


Figure 25. Correlation between moral identity and boundary separation of good condition; moral self-image and drift rate of good condition



*Figure 26.* Correlation between personal distance and boundary separation of neutral condition