

<sup>1</sup> Positivity bias in perceptual matching may reflect a spontaneous self-referential processing

<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> TBA

<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> Langenbeckstr. 1, Neuroimaging Center, University Medical Center Mainz, 55131 Mainz,

<sup>15</sup> Germany. E-mail: hcp4715@gmail.com

16

## Abstract

17 To navigate in a complex social world, individual has learnt to prioritize valuable  
18 information. Previous studies suggested the moral related stimuli was prioritized  
19 (Anderson, Siegel, Bliss-Moreau, & Barrett, 2011; Gantman & Van Bavel, 2014). Using  
20 social associative learning paradigm (self-tagging paradigm), we found that when geometric  
21 shapes, without soical meaning, were associated with different moral valence (morally  
22 good, neutral, or bad), the shapes that associated with positive moral valence were  
23 prioritized in a perceptual matching task. This patterns of results were robust across  
24 different procedures. Further, we tested whether this positivity effect was modulated by  
25 self-relevance by manipulating the self-relevance explicitly and found that this moral  
26 positivity effect was strong when the moral valence is describing oneself, but only weak  
27 evidence that such effect occured when the moral valence was describing others. We further  
28 found that this effect exist even when the self-relevance or the moral valence were  
29 presented as a task-irrelevant information, though the effect size become smaller. We also  
30 tested whether the positivity effect only exist in moral domain and found that this effect  
31 was not limited to moral domain. Exploratory analyses on task-questionnaire relationship  
32 found that moral self-image score (how closely one feel they are to the ideal moral image of  
33 themselves) is positively correlated to the  $d'$  of morally positive condition in singal  
34 detection and the drift rate using DDM, while the self-esteem is negatively correlated with  
35  $d'$  of neutral and morally negative conditions. These results suggest that the positive self  
36 prioritization in perceptual decision-making may reflect ...

37

*Keywords:* Perceptual decision-making, Self, positive bias, morality

38

Word count: X

39 Positivity bias in perceptual matching may reflect a spontaneous self-referential processing

40 # Introduction

41 Morality is the central of social life, regardless of the culture. It's so entrenched in the  
42 social life that one basic dimension of social perception is morality, i.e., whether a person is  
43 moral good or not. Also, morality is the most defining trait of a person. In the same vein,  
44 to succeed navigate in the social world, one has to consistently maintain a good moral  
45 reputation, i.e., moral self-image (to self), moral character (to others), because moral  
46 character of a person may determine whether she/he can find potential cooperators in the  
47 social world. This motivation may have strong influence on individual's behavioral and  
48 mind. For example, across three social domains, morality, social competence, and  
49 competence, people have the strongest self-enhancement effect in morality domain.

50 Previously studies revealed that people can maintain their moral self image even after their  
51 own unethical behaviors (e.g., cheating). Also, when asked how likely they will do ethical  
52 or unethical things, most participants showed the tendency of less likely to do unethical  
53 things []. Thus, the importance of morality in social life and most of adults are socialized  
54 individual, moral character is an dispensable part of their identity, the one that one need to  
55 maintain, even through self-deception.

56 From the information processing perspective, it is interesting, then, how moral  
57 self-image related information were processed by individuals? One might expect that the  
58 self-related moral information is more likely to be processed, elaborated, and perserved, so  
59 that the available information to the participant are more likely to consistent with one's  
60 expectation: good self! Indeed, recent studies found that memories retrieved are biased by  
61 their relevance to moral self-image and that this effect is modulated by the motivation of  
62 maintaining a positive moral self-image [].

63 It's less known, however, how the moral self-image related information is processed at  
64 the first lower cognitive stage, for example, in perceptual process where participant need to

65 decide their priority of information processing who the information should be processed.  
66 Though rarely appeared in moral related studies, the role of perception in understanding  
67 our cognition had been long root in psychology, social psychology is not exceptional. In  
68 1950s, Bruner (1957) had proposed the “New Look” approach of perception, which was  
69 resurrected by accumulating evidence (Stolier & Freeman, 2016; Xiao, Coppin, & Bavel,  
70 2016). These studies supported the view that there is a bidirectional interplay between  
71 perception and higher-level cognition, such as stereotype (Stolier & Freeman, 2016; Xiao et  
72 al., 2016), and self-relevance (Sui, He, & Humphreys, 2012). Few studies also tested  
73 whether moral-laden information was prioritized in the perception (Anderson et al., 2011;  
74 Gantman & Van Bavel, 2014). Still, the moral self-image information has rarely been  
75 studied (except Hu, Lan, Macrae, and Sui (2020))

76 These observation supported the view that the motivation of maintaining a positive  
77 self-image, especially on the aspect that are crucial to oneself [], can distort the cognitive  
78 processing [Greenwarld, 1988]. All the available model and evidence suggest that this  
79 positive self motivation also impact our perception processing.

80 To test this hypothesis, we adopted a self-tagging paradigm.

81 Outline:

82 Perceptual decision-making is an important window for understanding the cognition  
83 (Shadlen & Kiani, 2013), it's also has bi-directional interplay with higher-level cognition  
84 (Bruner, 1957; Gilbert & Li, 2013). Recently, there are accumulating evidence that  
85 perception also has bidirectional interplay with social cognition, such as morality  
86 (Anderson et al., 2011; Gantman & Van Bavel, 2014), stereotype (Stolier & Freeman, 2016;  
87 Xiao et al., 2016), and self-relevance (Sui et al., 2012).

88 One important feature of perception is the salience of the stimuli around us. Previous  
89 studies showed that emotional stimuli are salient and are prioritized. However, important  
90 modulate is whether the stimuli is self-relevant when they are physically and emotional

91 equally salient.

92 Here we investigated how the instantly learned moral valence changed the perceptual  
93 decision-making and the underlying psychological processes.

94 Given the importance of morality in social life (DeScioli, 2016) and identity (Freitas,  
95 Cikara, Grossmann, & Schlegel, 2017; Strohminger, Knobe, & Newman, 2017; Zhang,  
96 Chen, Schlegel, Hicks, & Chen, 2019), and evidence that moral character impacts how  
97 people evaluate themselves [XXXX], desired personality change (Sun & Goodwin, 2020),  
98 memory (Carlson, Maréchal, Oud, Fehr, & Crockett, 2020; Kouchaki & Gino, 2016; Shu,  
99 Gino, & Bazerman, 2011; Stanley & De Brigard, 2019), one would expect that moral  
100 character, the morality related trait, will also reflected in perceptual decision-making. Yet,  
101 this effect was less studied . This moral perception effect was driven by motivation  
102 (Gantman & Van Bavel, 2016) and influenced the information spreading online (Brady,  
103 Gantman, & Van Bavel, 2020). Especially lacking is the process of this effect.

104 The current study first explored and confirmed a positive effect of moral character in  
105 perceptual decision-making, using an associative learning task, then attempted to provide a  
106 mechanistic explanation for this positivity effect: spontaneous self-identification with the  
107 moral good character (Juechems and Summerfield (2019): Self-relevance is an important  
108 dimension in determine the value, i.e., intrinsic goal), both implicitly and explicitly.  
109 Finally, this positivity effect was also found in other social traits (beauty) but not  
110 non-social, emotional states.

111 Potential theoretical discussion points: Close distance of the semantic representation  
112 of self and moral character (attractor network) (Freeman & Ambady, 2011). The  
113 core/true/authentic self concept. social meter theory of self-esteem. evolutionary  
114 perspective of morality and moral self-conception, moral identity.

115 We reported behavioral results from eleven experiments. In first set of experiments,  
116 we found that shapes associated with morally positive person label were responded faster

and more accurately. In the second set of experiments, we explore the potential role of good self in perceptual matching task and added one more independent variable, we found that the effect was mainly on good self. In the third part we tested whether the morality will automatically binds with self but not other. Finally, we explore the correlation between behavioral task and questionnaire scores.

122

## Disclosures

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

127 All the experiments reported were not pre-registered. Most experiments (1a ~ 6b,  
128 except experiment 3b) reported in the current study were first finished between 2014 to  
129 2016 in Tsinghua University, Beijing, China. Participants in these experiments were  
130 recruited in the local community. To increase the sample size of experiments to 50 or more  
131 (Simmons, Nelson, & Simonsohn, 2013), we recruited additional participants in Wenzhou  
132 University, Wenzhou, China in 2017 for experiment 1a, 1b, 4a, and 4b. Experiment 3b was  
133 finished in Wenzhou University in 2017. To have a better estimation of the effect size, we  
134 included the data from two experiments (experiment 7a, 7b) that were reported in Hu et  
135 al. (2020) (See Table S1 for overview of these experiments).

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

**General methods****139 Design and Procedure**

140 This series of experiments started to test the effect of instantly acquired true self  
141 (moral self) on perceptual decision-making. For this purpose, we used the social associative  
142 learning paradigm (or tagging paradigm)(Sui et al., 2012), in which participants first  
143 learned the associations between geometric shapes and labels of person with different moral  
144 character (e.g., in first three studies, the triangle, square, and circle and good person,  
145 neutral person, and bad person, respectively). The associations of the shapes and label  
146 were counterbalanced across participants. After remembered the associations, participants  
147 finished a practice phase to familiar with the task, in which they viewed one of the shapes  
148 upon the fixation while one of the labels below the fixation and judged whether the shape  
149 and the label matched the association they learned. When participants reached 60% or  
150 higher accuracy at the end of the practicing session, they started the experimental task  
151 which was the same as in the practice phase.

152 The experiment 1a, 1b, 1c, 2, and 6a shared a 2 (matching: match vs. nonmatch) by  
153 3 (moral valence: good vs. neutral vs. bad) within-subject design. Experiment 1a was the  
154 first one of the whole series studies and 1b, 1c, and 2 were conducted to exclude the  
155 potential confounding factors. More specifically, experiment 1b used different Chinese  
156 words as label to test whether the effect only occurred with certain familiar words.  
157 Experiment 1c manipulated the moral valence indirectly: participants first learned to  
158 associate different moral behaviors with different neutral names, after remembered the  
159 association, they then performed the perceptual matching task by associating names with  
160 different shapes. Experiment 2 further tested whether the way we presented the stimuli  
161 influence the effect of valence, by sequentially presenting labels and shapes. Note that part  
162 of participants of experiment 2 were from experiment 1a because we originally planned a  
163 cross task comparison. Experiment 6a, which shared the same design as experiment 2, was

<sup>164</sup> an EEG experiment which aimed at exploring the neural correlates of the effect. But we  
<sup>165</sup> will focus on the behavioral results of experiment 6a in the current manuscript.

<sup>166</sup> For experiment 3a, 3b, 4a, 4b, 6b, 7a, and 7b, we included self-reference as another  
<sup>167</sup> within-subject variable in the experimental design. For example, the experiment 3a directly  
<sup>168</sup> extend the design of experiment 1a into a 2 (matchness: match vs. nonmatch) by 2  
<sup>169</sup> (reference: self vs. other) by 3 (moral valence: good vs. neutral vs. bad) within-subject  
<sup>170</sup> design. Thus in experiment 3a, there were six conditions (good-self, neutral-self, bad-self,  
<sup>171</sup> good-other, neutral-other, and bad-other) and six shapes (triangle, square, circle, diamond,  
<sup>172</sup> pentagon, and trapezoids). The experiment 6b was an EEG experiment extended from  
<sup>173</sup> experiment 3a but presented the label and shape sequentially. Because of the relatively  
<sup>174</sup> high working memory load (six label-shape pairs), experiment 6b were conducted in two  
<sup>175</sup> days: the first day participants finished perceptual matching task as a practice, and the  
<sup>176</sup> second day, they finished the task again while the EEG signals were recorded. Experiment  
<sup>177</sup> 3b was designed to separate the self-referential trials and other-referential trials. That is,  
<sup>178</sup> participants finished two different blocks: in the self-referential blocks, they only responded  
<sup>179</sup> to good-self, neutral-self, and bad-self, with half match trials and half non-match trials; for  
<sup>180</sup> the other-reference blocks, they only responded to good-other, neutral-other, and  
<sup>181</sup> bad-other. Experiment 7a and 7b were designed to test the cross task robustness of the  
<sup>182</sup> effect we observed in the aforementioned experiments (see, Hu et al., 2020). The matching  
<sup>183</sup> task in these two experiments shared the same design with experiment 3a, but only with  
<sup>184</sup> two moral valence, i.e., good vs. bad. We didn't include the neutral condition in  
<sup>185</sup> experiment 7a and 7b because we found that the neutral and bad conditions constantly  
<sup>186</sup> showed non-significant results in experiment 1 ~ 6.

<sup>187</sup> Experiment 4a and 4b were design to test the automaticity of the binding between  
<sup>188</sup> self/other and moral valence. In 4a, we used only two labels (self vs. other) and two shapes  
<sup>189</sup> (circle, square). To manipulate the moral valence, we added the moral-related words within  
<sup>190</sup> the shape and instructed participants to ignore the words in the shape during the task. In

191 4b, we reversed the role of self-reference and valence in the task: participant learnt three  
192 labels (good-person, neutral-person, and bad-person) and three shapes (circle, square, and  
193 triangle), and the words related to identity, “self” or “other”, were presented in the shapes.  
194 As in 4a, participants were told to ignore the words inside the shape during the task.

195 Finally, experiment 5 was design to test the specificity of the moral valence. We  
196 extended experiment 1a with an additional independent variable: domains of the valence  
197 words. More specifically, besides the moral valence, we also added valence from other  
198 domains: appearance of person (beautiful, neutral, ugly), appearance of a scene (beautiful,  
199 neutral, ugly), and emotion (happy, neutral, and sad). Label-shape pairs from different  
200 domains were separated into different blocks.

201 E-prime 2.0 was used for presenting stimuli and collecting behavioral responses,  
202 except that experiment 7a and 7b used Matlab Psychtoolbox (Brainard, 1997; Pelli, 1997).  
203 For participants recruited in Tsinghua University, they finished the experiment individually  
204 in a dim-lighted chamber, stimuli were presented on 22-inch CRT monitors and their head  
205 were fixed by a chin-rest brace. The distance between participants’ eyes and the screen was  
206 about 60 cm. The visual angle of geometric shapes was about  $3.7^\circ \times 3.7^\circ$ , the fixation cross  
207 is of ( $0.8^\circ \times 0.8^\circ$  of visual angle) at the center of the screen. The words were of  $3.6^\circ \times 1.6^\circ$   
208 visual angle. The distance between the center of the shape or the word and the fixation  
209 cross was  $3.5^\circ$  of visual angle. For participants recruited in Wenzhou University, they  
210 finished the experiment in a group consisted of 3 ~ 12 participants in a dim-lighted testing  
211 room. Participants were required to finished the whole experiment independently. Also,  
212 they were instructed to start the experiment at the same time, so that the distraction  
213 between participants were minimized. The stimuli were presented on 19-inch CRT monitor.  
214 The visual angles are could not be exactly controlled because participants’s chin were not  
215 fixed.

216 In most of these experiments, participant were also asked to fill a battery of

<sup>217</sup> questionnaire after they finish the behavioral tasks. All the questionnaire data are open  
<sup>218</sup> (see, dataset 4 in Liu et al., 2020). See Table S1 for a summary information about all the  
<sup>219</sup> experiments.

<sup>220</sup> **Data analysis**

<sup>221</sup> **Analysis of individual study.** We used the `tidyverse` of r (see script  
<sup>222</sup> `Load_save_data.r`) to exclude the practicing trials, invalid trials of each participants, and  
<sup>223</sup> invalid participants, if there were any, in the raw data. Results of each experiment were  
<sup>224</sup> then analyzed in three different approaches.

<sup>225</sup> ***Classic NHST.***

<sup>226</sup> First, as in Sui et al. (2012), we analyzed the accuracy and reaction times using  
<sup>227</sup> classic repeated measures ANOVA in the Null Hypothesis Significance Test (NHST)  
<sup>228</sup> framework. Repeated measures ANOVAs is essentially a two-step mixed model. In the first  
<sup>229</sup> step, we estimate the parameter on individual level, and in the second step, we used  
<sup>230</sup> repeated ANOVA to test the Null hypothesis. More specifically, for the accuracy, we used a  
<sup>231</sup> signal detection approach, in which individual' sensitivity  $d'$  was estimated first. To  
<sup>232</sup> estimate the sensitivity, we treated the match condition as the signal while the nonmatch  
<sup>233</sup> conditions as noise. Trials without response were coded either as “miss” (match trials) or  
<sup>234</sup> “false alarm” (nonmatch trials). Given that the match and nonmatch trials are presented  
<sup>235</sup> in the same way and had same number of trials across all studies, we assume that  
<sup>236</sup> participants' inner distribution of these two types of trials had equal variance but may had  
<sup>237</sup> different means. That is, we used the equal variance Gaussian SDT model (EVSDT) here  
<sup>238</sup> (Rouder & Lu, 2005). The  $d'$  was then estimated as the difference of the standardized hit  
<sup>239</sup> and false alarm rats (Stanislaw & Todorov, 1999):

$$d' = zHR - zFAR = \Phi^{-1}(HR) - \Phi^{-1}(FAR)$$

<sup>240</sup> where the  $HR$  means hit rate and the  $FAR$  mean false alarm rate.  $zHR$  and  $zFAR$  are

241 the standardized hit rate and false alarm rates, respectively. These two  $z$ -scores were  
 242 converted from proportion (i.e., hit rate or false alarm rate) by inverse cumulative normal  
 243 density function,  $\Phi^{-1}$  ( $\Phi$  is the cumulative normal density function, and is used convert  $z$   
 244 score into probabilities). Another parameter of signal detection theory, response criterion  $c$ ,  
 245 is defined by the negative standardized false alarm rate (DeCarlo, 1998):  $-zFAR$ .

246 For the reaction times (RTs), only RTs of accurate trials were analyzed. We first  
 247 calculate the mean RTs of each participant and then subject the mean RTs of each  
 248 participant to repeated measures ANOVA. Note that we set the alpha as .05. The repeated  
 249 measure ANOVA was done by `afex` package (<https://github.com/singmann/afex>).

250 To control the false positive rate when conducting the post-hoc comparisons, we used  
 251 Bonferroni correction.

252 ***Bayesian hierarchical generalized linear model (GLM).***

253 The classic NHST approach may ignore the uncertainty in estimate of the parameters  
 254 for SDT (Rouder & Lu, 2005), and using mean RT assumes normal distribution of RT  
 255 data, which is always not true because RTs distribution is skewed (Rousselet & Wilcox,  
 256 2019). To better estimate the uncertainty and use a more appropriate model, we also tried  
 257 Bayesian hierarchical generalized linear model to analyze each experiment's accuracy and  
 258 RTs data. We used BRMs (Bürkner, 2017) to build the model, which used Stan (Carpenter  
 259 et al., 2017) to estimate the posterior.

260 In the GLM model, we assume that the accuracy of each trial is Bernoulli distributed  
 261 (binomial with 1 trial), with probability  $p_i$  that  $y_i = 1$ .

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

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

$$\Phi(p_i) = \beta_0 + \beta_1 IsMatch_i * Valence_i$$

264 The outcomes  $y_i$  are 0 if the participant responded “nonmatch” on trial  $i$ , 1 if they  
 265 responded “match”. The probability of the “match” response for trial  $i$  for a participant is  
 266  $p_i$ . We then write the generalized linear model on the probits (z-scores;  $\Phi$ , “Phi”) of  $ps$ .  $\Phi$   
 267 is the cumulative normal density function and maps  $z$  scores to probabilities. Given this  
 268 parameterization, the intercept of the model ( $\beta_0$ ) is the standardized false alarm rate  
 269 (probability of saying 1 when predictor is 0), which we take as our criterion  $c$ . The slope of  
 270 the model ( $\beta_1$ ) is the increase of saying 1 when predictor is 1, in  $z$ -scores, which is another  
 271 expression of  $d'$ . Therefore,  $c = -zHR = -\beta_0$ , and  $d' = \beta_1$ .

272 In each experiment, we had multiple participants, then we need also consider the  
 273 variations between subjects, i.e., a hierarchical mode in which individual’s parameter and  
 274 the population level parameter are estimated simultaneously. We assume that the  
 275 outcome of each trial is Bernoulli distributed (binomial with 1 trial), with probability  $p_{ij}$   
 276 that  $y_{ij} = 1$ .

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

277 Similarly, the generalized linear model was extended to two levels:

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

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

281 The subjective-specific intercepts ( $\beta_0 = -zFAR$ ) and slopes ( $\beta_1 = d'$ ) are describe  
 282 by multivariate normal with means and a covariance matrix for the parameters.

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

283 For the reaction time, we used the log normal distribution

284 ([https://lindeloev.github.io/shiny-rt/#34\\_\(shifted\)\\_log-normal](https://lindeloev.github.io/shiny-rt/#34_(shifted)_log-normal)). This distribution has  
 285 two parameters:  $\mu$ ,  $\sigma$ .  $\mu$  is the mean of the logNormal distribution, and  $\sigma$  is the disperse of  
 286 the distribution. The log normal distribution can be extended to shifted log normal  
 287 distribution, with one more parameter: shift, which is the earliest possible response.

$$y_i = \beta_0 + \beta_1 * IsMatch_i * Valence_i$$

288 Shifted log-normal distribution:

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

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

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

$$\sigma_j \sim Cauchy()$$

### 290 ***Hierarchical drift diffusion model (HDDM).***

291 To further explore the psychological mechanism under perceptual decision-making, we

292 used HDDM (Wiecki, Sofer, & Frank, 2013) to model our RTs and accuracy data. We used  
 293 the prior implemented in HDDM, that is, informative priors that constrains parameter  
 294 estimates to be in the range of plausible values based on past literature (Matzke &  
 295 Wagenmakers, 2009). As reported in Hu et al. (2020), we used the response code approach,  
 296 match response were coded as 1 and nonmatch responses were coded as 0. To fully explore  
 297 all parameters, we allow all four parameters of DDM free to vary. We then extracted the  
 298 estimation of all the four parameters for each participants for the correlation analyses.

299 However, because the starting point is only related to response (match vs. non-match) but  
 300 not the valence of the stimuli, we didn't included it in correlation analysis.

301 **Synthesized results.** We also reported the synthesized results from the  
302 experiments, because many of them shared the similar experimental design. We reported  
303 the results in five parts: valence effect, explicit interaction between valence and  
304 self-relevance, implicit interaction between valence and self-relevance, specificity of valence  
305 effect, and behavior-questionnaire correlation.

306 For the first two parts, we reported the synthesized results from Frequentist's  
307 approach(mini-meta-analysis, Goh, Hall, & Rosenthal, 2016). The mini meta-analyses were  
308 carried out by using `metafor` package (Viechtbauer, 2010). We first calculated the mean of  
309  $d'$  and RT of each condition for each participant, then calculate the effect size (Cohen's  $d$ )  
310 and variance of the effect size for all contrast we interested: Good v. Bad, Good v.  
311 Neutral, and Bad v. Neutral for the effect of valence, and self vs. other for the effect of  
312 self-relevance. Cohen's  $d$  and its variance were estimated using the following formula  
313 (Cooper, Hedges, & Valentine, 2009):

$$d = \frac{(M_1 - M_2)}{\sqrt{(sd_1^2 + sd_2^2) - 2rsd_1sd_2}}\sqrt{2(1-r)}$$

$$var.d = 2(1-r)\left(\frac{1}{n} + \frac{d^2}{2n}\right)$$

314  $M_1$  is the mean of the first condition,  $sd_1$  is the standard deviation of the first  
315 condition, while  $M_2$  is the mean of the second condition,  $sd_2$  is the standard deviation of  
316 the second condition.  $r$  is the correlation coefficient between data from first and second  
317 condition.  $n$  is the number of data point (in our case the number of participants included  
318 in our research).

319 The effect size from each experiment were then synthesized by random effect model  
320 using `metafor` (Viechtbauer, 2010). Note that to avoid the cases that some participants  
321 participated more than one experiments, we inspected the all available information of  
322 participants and only included participants' results from their first participation. As

<sup>323</sup> mentioned above, 24 participants were intentionally recruited to participate both exp 1a  
<sup>324</sup> and exp 2, we only included their results from experiment 1a in the meta-analysis.

<sup>325</sup> We also estimated the synthesized effect size using Bayesian hierarchical model,  
<sup>326</sup> which extended the two-level hierarchical model in each experiment into three-level model,  
<sup>327</sup> which experiment as an additional level. For SDT, we can use a nested hierarchical model  
<sup>328</sup> to model all the experiment with similar design:

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

<sup>329</sup> where

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

<sup>330</sup> The outcomes  $y_{ijk}$  are 0 if participant  $j$  in experiment k responded “nonmatch” on trial  $i$ ,  
<sup>331</sup> 1 if they responded “match”.

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

<sup>332</sup> and the experiment level parameter  $mu_{0k}$  and  $mu_{1k}$  is from a higher order  
<sup>333</sup> distribution:

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

<sup>334</sup> in which  $\mu_0$  and  $\mu_1$  means the population level parameter.

<sup>335</sup> This model can be easily expand to three-level model in which participants and  
<sup>336</sup> experiments are two group level variable and participants were nested in the experiments.

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

<sup>337</sup>  $y_{ijk}$  is the RT of the  $i$ th trial of the  $j$ th participants in the  $k$ th experiment.

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

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

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

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

338        Using the Bayesian hierarchical model, we can directly estimate the over-all effect of  
 339        valence on  $d'$  across all experiments with similar experimental design, instead of using a  
 340        two-step approach where we first estimate the  $d'$  for each participant and then use a  
 341        random effect model meta-analysis (Goh et al., 2016).

342        ***Valence effect.***

343        We synthesized effect size of  $d'$  and RT from experiment 1a, 1b, 1c, 2, 5 and 6a for  
 344        the valence effect. We reported the synthesized the effect across all experiments that tested  
 345        the valence effect, using the mini meta-analysis approach (Goh et al., 2016).

346        ***Explicit interaction between Valence and self-relevance.***

347        The results from experiment 3a, 3b, 6b, 7a, and 7b. These experiments explicitly  
 348        included both moral valence and self-reference.

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

350        In the third part, we focused on experiment 4a and 4b, which were designed to  
 351        examine the implicit effect of the interaction between moral valence and self-referential  
 352        processing. We are interested in one particular question: will self-referential and morally  
 353        positive valence had a mutual facilitation effect. That is, when moral valence (experiment

354 4a) or self-referential (experiment 4a) was presented as task-irrelevant stimuli, whether  
355 they would facilitate self-referential or valence effect on perceptual decision-making. For  
356 experiment 4a, we reported the comparisons between different valence conditions under the  
357 self-referential task and other-referential task. For experiment 4b, we first calculated the  
358 effect of valence for both self- and other-referential conditions and then compared the effect  
359 size of these three contrast from self-referential condition and from other-referential  
360 condition. Note that the results were also analyzed in a standard repeated measure  
361 ANOVA (see supplementary materials).

362 ***Specificity of the valence effect.***

363 In this part, we reported the data from experiment 5, which included positive,  
364 neutral, and negative valence from four different domains: morality, aesthetic of person,  
365 aesthetic of scene, and emotion. This experiment was design to test whether the positive  
366 bias is specific to morality.

367 ***Behavior-Questionnaire correlation.***

368 Finally, we explored correlation between results from behavioral results and  
369 self-reported measures.

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

375 For the behavioral task part, we used three parameters from drift diffusion model:  
376 drift rate ( $v$ ), boundary separation ( $a$ ), and non decision-making time ( $t$ ), because these  
377 parameters has relative clear psychological meaning. We used the mean of parameter  
378 posterior distribution as the estimate of each parameter for each participants in the  
379 correlation analysis.

380 Based on results from the experiment, we reason that the correlation between  
381 behavioral result in self-referential will appear in the data without mentioning the  
382 self/other (exp 3a, 3b, 6a, 7a, 7b). To this end, we first calculated the correlation between  
383 behavioral indicators and questionnaires for self-referential and other-referential separately.  
384 Given the small sample size of the data ( $N =$ ), we used a relative liberal threshold for  
385 these explorations ( $\alpha = 0.1$ ).

386 Then we confirmed the significant results from the data without self- and  
387 other-referential (exp1a, 1b, 1c, 2, 5, 6a). In this confirmatory analysis, we used  $\alpha =$   
388 0.05 and used bootstrap by BootES package (Kirby & Gerlanc, 2013) to estimate the  
389 correlation. To avoid false positive, we further determined the threshold for significant by  
390 permutation. More specifically, for each pairs that initially with  $p < .05$ , we randomly  
391 shuffle the participants data of each score and calculated the correlation between the  
392 shuffled vectors. After repeating this procedure for 5000 times, we choose arrange these  
393 5000 correlation coefficients and use the 95% percentile number as our threshold.

### 394 **Part 1: Moral valence effect**

395 In this part, we report five experiments that aimed at testing whether the instantly  
396 acquired association between shapes and good person would be prioritized in perceptual  
397 decision-making.

#### 398 **Experiment 1a**

##### 399 **Methods.**

##### 400 ***Participants.***

401 57 college students (38 female, age =  $20.75 \pm 2.54$  years) participated. 39 of them  
402 were recruited from Tsinghua University community in 2014; 18 were recruited from  
403 Wenzhou University in 2017. All participants were right-handed except one, and all had

404 normal or corrected-to-normal vision. Informed consent was obtained from all participants  
405 prior to the experiment according to procedures approved by the local ethics committees. 6  
406 participant's data were excluded from analysis because nearly random level of accuracy,  
407 leaving 51 participants (34 female, age =  $20.72 \pm 2.44$  years).

408        ***Stimuli and Tasks.***

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

412        ***Procedure.***

413        This experiment had two phases. First, there was a brief learning stage. Participants  
414 were asked to learn the relationship between geometric shapes (triangle, square, and circle)  
415 and different person (bad person, a good person, or a neutral person). For example, a  
416 participant was told, "bad person is a circle; good person is a triangle; and a neutral person  
417 is represented by a square." After participant remember the associations (usually in a few  
418 minutes), participants started a practicing phase of matching task which has the exact task  
419 as in the experimental task. In the experimental task, participants judged whether  
420 shape-label pairs, which were subsequently presented, were correct. Each trial started with  
421 the presentation of a central fixation cross for 500 ms. Subsequently, a pairing of a shape  
422 and label (good person, bad person, and neutral person) was presented for 100 ms. The  
423 pair presented could confirm to the verbal instruction for each pairing given in the training  
424 stage, or it could be a recombination of a shape with a different label, with the shape-label  
425 pairings being generated at random. The next frame showed a blank for 1100ms.  
426 Participants were expected to judge whether the shape was correctly assigned to the person  
427 by pressing one of the two response buttons as quickly and accurately as possible within  
428 this timeframe (to encourage immediate responding). Feedback (correct or incorrect) was  
429 given on the screen for 500 ms at the end of each trial, if no response detected, "too slow"

430 was presented to remind participants to accelerate. Participants were informed of their  
431 overall accuracy at the end of each block. The practice phase finished and the experimental  
432 task began after the overall performance of accuracy during practice phase achieved 60%.  
433 For participants from the Tsinghua community, they completed 6 experimental blocks of 60  
434 trials. Thus, there were 60 trials in each condition (bad-person match, bad-person  
435 nonmatch, good-person match, good-person nonmatch, neutral-person match, and  
436 neutral-person nonmatch). For the participants from Wenzhou University, they finished 6  
437 blocks of 120 trials, therefore, 120 trials for each condition.

438 ***Data analysis.***

439 As described in general methods section, this experiment used three approaches to  
440 analyze the behavioral data: Classical NHST, Bayesian Hierarchical Generalized Linear  
441 Model, and Hierarchical drift diffusion model.

442 **Results.**

443 ***Classic NHST.***

444 *d prime.*

445 Figure 1 shows *d prime* and reaction times during the perceptual matching task. We  
446 conducted a single factor (valence: good, neutral, bad) repeated measure ANOVA.

447 We found the effect of Valence ( $F(1.96, 97.84) = 6.19$ ,  $MSE = 0.27$ ,  $p = .003$ ,  
448  $\hat{\eta}_G^2 = .020$ ). The post-hoc comparison with multiple comparison correction revealed that  
449 the shapes associated with Good-person (2.11, SE = 0.14) has greater *d prime* than shapes  
450 associated with Bad-person (1.75, SE = 0.14),  $t(50) = 3.304$ ,  $p = 0.0049$ . The Good-person  
451 condition was also greater than the Neutral-person condition (1.95, SE = 0.16), but didn't  
452 reach statistical significant,  $t(50) = 1.54$ ,  $p = 0.28$ . Neither the Neutral-person condition is  
453 significantly greater than the Bad-person condition,  $t(50) = 2.109$ ,  $p = .098$ .

454 ***Reaction times.***

455 We conducted 2 (Matchness: match v. nonmatch) by 3 (Valence: good, neutral, bad)

456 repeated measure ANOVA. We found the main effect of Matchness ( $F(1, 50) = 232.39$ ,

457  $MSE = 948.92, p < .001, \hat{\eta}_G^2 = .104$ ), main effect of valence ( $F(1.87, 93.31) = 9.62$ ,

458  $MSE = 1,673.86, p < .001, \hat{\eta}_G^2 = .016$ ), and interaction between Matchness and Valence

459 ( $F(1.73, 86.65) = 8.52, MSE = 1,441.75, p = .001, \hat{\eta}_G^2 = .011$ ).

460 We then carried out two separate ANOVA for Match and Mismatched trials. For

461 matched trials, we found the effect of valence . We further examined the effect of valence

462 for both self and other for matched trials. We found that shapes associated with Good

463 Person (684 ms, SE = 11.5) responded faster than Neutral (709 ms, SE = 11.5),  $t(50) =$

464 -2.265,  $p = 0.0702$ ) and Bad Person (728 ms, SE = 11.7),  $t(50) = -4.41, p = 0.0002$ ), and

465 the Neutral condition was faster than the Bad condition,  $t(50) = -2.495, p = 0.0415$ ). For

466 non-matched trials, there was no significant effect of Valence ()�.

### 467 ***Bayesian hierarchical GLM.***

468  $d'$  prime.

469 We fitted a Bayesian hierarchical GLM for signal detection theory approach. The

470 results showed that when the shapes were tagged with labels with different moral valence,

471 the sensitivity ( $d'$ ) and criteria ( $c$ ) were both influence. For the  $d'$ , we found that the

472 shapes tagged with morally good person (2.46, 95% CI[2.21 2.72]) is greater than shapes

473 tagged with moral bad (2.07, 95% CI[1.83 2.32]),  $P_{PosteriorComparison} = 1$ . Shape tagged

474 with morally good person is also greater than shapes tagged with neutral person (2.23,

475 95% CI[1.95 2.49]),  $P_{PosteriorComparison} = 0.97$ . Also, the shapes tagged with neutral

476 person is greater than shapes tagged with morally bad person,  $P_{PosteriorComparison} = 0.92$ .

477 Interesting, we also found the criteria for three conditions also differ, the shapes

478 tagged with good person has the highest criteria (-1.01, [-1.14 -0.88]), followed by shapes

479 tagged with neutral person(-1.06, [-1.21 -0.92]), and then the shapes tagged with bad

480 person(-1.11, [-1.25 -0.97]). However, pair-wise comparison showed that only showed strong

481 evidence for the difference between good and bad conditions.

482 *Reaction times.*

483 We fitted a Bayesian hierarchical GLM for RTs, with a log-normal distribution as the  
484 link function. We used the posterior distribution of the regression coefficient to make  
485 statistical inferences. As in previous studies, the matched conditions are much faster than  
486 the mismatched trials ( $P_{PosteriorComparison} = 1$ ). We focused on matched trials only, and  
487 compared different conditions: Good is faster than the neutral,  $P_{PosteriorComparison} = .99$ ,  
488 it was also faster than the Bad condition,  $P_{PosteriorComparison} = 1$ . And the neutral  
489 condition is faster than the bad condition,  $P_{PosteriorComparison} = .99$ . However, the  
490 mismatched trials are largely overlapped. See Figure 2.

491 **HDDM.**

492 We fitted our data with HDDM, using the response-coding (See also, Hu et al., 2020).  
493 We estimated separate drift rate ( $v$ ), non-decision time ( $T_0$ ), and boundary separation ( $a$ )  
494 for each condition. We found that the shapes tagged with good person has higher drift rate  
495 and higher boundary separation than shapes tagged with both neutral and bad person.  
496 Also, the shapes tagged with neutral person has a higher drift rate than shapes tagged  
497 with bad person, but not for the boundary separation. Finally, we found that shapes  
498 tagged with bad person had longer non-decision time (see Figure 3).

499 **Experiment 1b**

500 In this study, we aimed at excluding the potential confounding factor of the  
501 familiarity of words we used in experiment 1a, by matching the familiarity of the words.

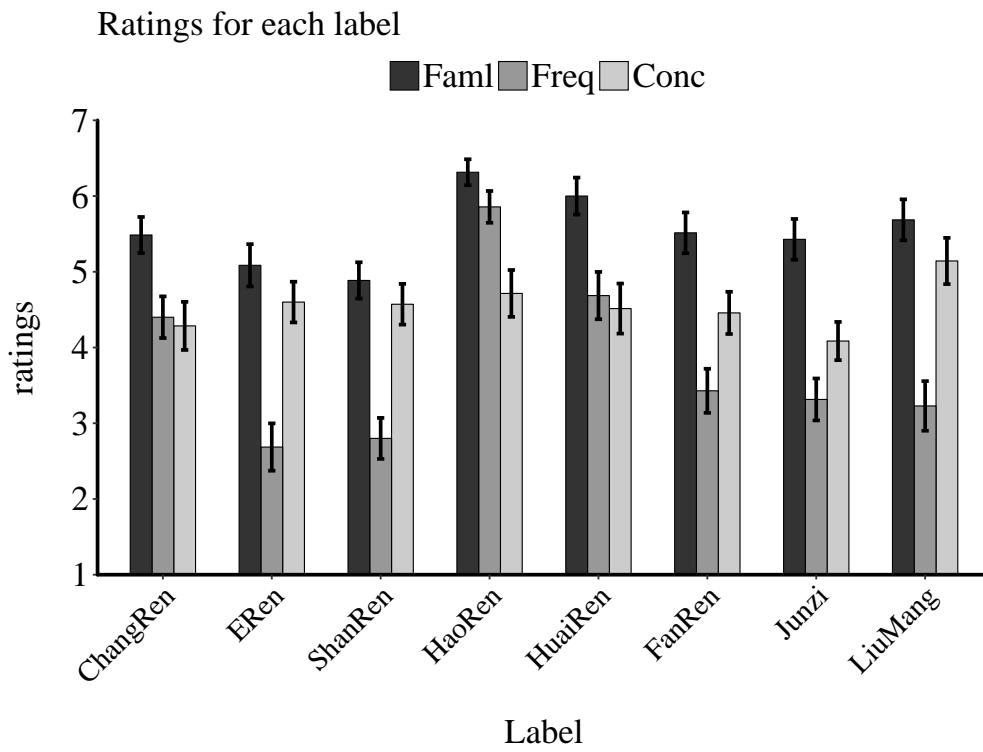
502 **Method.**

503 *Participants.*

504 72 college students (49 female, age =  $20.17 \pm 2.08$  years) participated. 39 of them  
505 were recruited from Tsinghua University community in 2014; 33 were recruited from

506 Wenzhou University in 2017. All participants were right-handed except one, and all had  
507 normal or corrected-to-normal vision. Informed consent was obtained from all participants  
508 prior to the experiment according to procedures approved by the local ethics committees.  
509 20 participant's data were excluded from analysis because nearly random level of accuracy,  
510 leaving 52 participants (36 female, age =  $20.25 \pm 2.31$  years).

511 **Stimuli and Tasks.** Three geometric shapes (triangle, square, and circle, with  $3.7^\circ$   
512  $\times 3.7^\circ$  of visual angle) were presented above a white fixation cross subtending  $0.8^\circ \times 0.8^\circ$   
513 of visual angle at the center of the screen. The three shapes were randomly assigned to  
514 three labels with different moral valence: a morally bad person (" ", ERen), a morally  
515 good person (" ", ShanRen) or a morally neutral person (" ", ChangRen). The order of  
516 the associations between shapes and labels was counterbalanced across participants. Three  
517 labels used in this experiment is selected based on the rating results from an independent  
518 survey, in which participants rated the familiarity, frequency, and concreteness of eight  
519 different words online. Of the eight words, three of them are morally positive (HaoRen,  
520 ShanRen, Junzi), two of them are morally neutral (ChangRen, FanRen), and three of them  
521 are morally negative (HuaiRen, ERen, LiuMang). An independent sample consist of 35  
522 participants (22 females, age  $20.6 \pm 3.11$ ) were recruited to rate these words. Based on the  
523 ratings (see supplementary materials Figure S1), we selected ShanRen, ChangRen, and  
524 ERen to represent morally positive, neutral, and negative person.



525

### Procedure.

526

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

529

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

530

### Results.

531

#### NHST.

532

Figure 4 shows  $d$  prime and reaction times of experiment 1b.

533

$d$  prime.

534

Repeated measures ANOVA revealed main effect of valence,  $F(1.83, 93.20) = 14.98$ ,

$MSE = 0.18$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .053$ . Paired t test showed that the Good-Person condition

$(1.87 \pm 0.102)$  was with greater  $d$  prime than Neutral condition  $(1.44 \pm 0.101$ ,  $t(51) =$

$5.945$ ,  $p < 0.001$ ). We also found that the Bad-Person condition  $(1.67 \pm 0.11)$  has also

$538$  greater  $d$  prime than neutral condition ,  $t(51) = 3.132$ ,  $p = 0.008$ ). There Good-person

539 condition was also slightly greater than the bad condition,  $t(51) = 2.265, p = 0.0701$ .

540 *Reaction times.*

541 We found interaction between Matchness and Valence ( $F(1.95, 99.31) = 19.71$ ,  
 542  $MSE = 960.92, p < .001, \hat{\eta}_G^2 = .031$ ) and then analyzed the matched trials and  
 543 mismatched trials separately, as in experiment 1a. For matched trials, we found the effect  
 544 of valence  $F(1.94, 99.10) = 33.97, MSE = 1,343.19, p < .001, \hat{\eta}_G^2 = .115$ . Post-hoc  $t$ -tests  
 545 revealed that shapes associated with Good Person ( $684 \pm 8.77$ ) were responded faster than  
 546 Neutral-Person ( $740 \pm 9.84$ ), ( $t(51) = -8.167, p < 0.001$ ) and Bad Person ( $728 \pm 9.15$ ),  
 547  $t(51) = -5.724, p < 0.0001$ ). While there was no significant differences between Neutral and  
 548 Bad-Person condition ( $t(51) = 1.686, p = 0.221$ ). For non-matched trials, there was no  
 549 significant effect of Valence ( $F(1.90, 97.13) = 1.80, MSE = 430.15, p = .173, \hat{\eta}_G^2 = .003$ ).

550 **BGLM.**

551 *Signal detection theory analysis of accuracy.*

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

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

565        *Reaction time.*

566        We fitted a Bayesian hierarchical GLM for RTs, with a log-normal distribution as the  
567        link function. We used the posterior distribution of the regression coefficient to make  
568        statistical inferences. As in previous studies, the matched conditions are much faster than  
569        the mismatched trials ( $P_{PosteriorComparison} = 1$ ). We focused on matched trials only, and  
570        compared different conditions: Good is faster than the neutral,  $P_{PosteriorComparison} = .99$ ,  
571        it was also faster than the Bad condition,  $P_{PosteriorComparison} = 1$ . And the neutral  
572        condition is faster than the bad condition,  $P_{PosteriorComparison} = .99$ . However, the  
573        mismatched trials are largely overlapped. See Figure 5.

574        **HDDM.**

575        We found that the shapes tagged with good person has higher drift rate and higher  
576        boundary separation than shapes tagged with both neutral and bad person. Also, the  
577        shapes tagged with neutral person has a higher drift rate than shapes tagged with bad  
578        person, but not for the boundary separation. Finally, we found that shapes tagged with  
579        bad person had longer non-decision time (see figure 6).

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

585        **Experiment 1c**

586        In this study, we further control the valence of words using in our experiment.  
587        Instead of using label with moral valence, we used valence-neutral names in China.  
588        Participant first learn behaviors of the different person, then, they associate the names and  
589        shapes. And then they perform a name-shape matching task.

590       **Method.**

591       ***Participants.***

592       23 college students (15 female, age =  $22.61 \pm 2.62$  years) participated. All of them  
593       were recruited from Tsinghua University community in 2014. Informed consent was  
594       obtained from all participants prior to the experiment according to procedures approved by  
595       the local ethics committees. No participant was excluded because they overall accuracy  
596       were above 0.6.

597       ***Stimuli and Tasks.***

598       Three geometric shapes (triangle, square, and circle, with  $3.7^\circ \times 3.7^\circ$  of visual angle)  
599       were presented above a white fixation cross subtending  $0.8^\circ \times 0.8^\circ$  of visual angle at the  
600       center of the screen. The three most common names were chosen, which are neutral in  
601       moral valence before the manipulation. Three names (Zhang, Wang, Li) were first paired  
602       with three paragraphs of behavioral description. Each description includes one sentence of  
603       biographic information and four sentences that describing the moral behavioral under that  
604       name. To assess the that these three descriptions represented good, neutral, and bad  
605       valence, we collected the ratings of three person on six dimensions: morality, likability,  
606       trustworthiness, dominance, competence, and aggressiveness, from an independent sample  
607       ( $n = 34$ , 18 female, age =  $19.6 \pm 2.05$ ). The rating results showed that the person with  
608       morally good behavioral description has higher score on morality ( $M = 3.59$ ,  $SD = 0.66$ )  
609       than neutral ( $M = 0.88$ ,  $SD = 1.1$ ),  $t(33) = 12.94$ ,  $p < .001$ , and bad conditions ( $M = -3.4$ ,  
610        $SD = 1.1$ ),  $t(33) = 30.78$ ,  $p < .001$ . Neutral condition was also significant higher than bad  
611       conditions  $t(33) = 13.9$ ,  $p < .001$  (See supplementary materials).

612       ***Procedure.***

613       After arriving the lab, participants were informed to complete two experimental  
614       tasks, first a social memory task to remember three person and their behaviors, after tested  
615       for their memory, they will finish a perceptual matching task. In the social memory task,

the descriptions of three person were presented without time limitation. Participant self-paced to memorized the behaviors of each person. After they memorizing, a recognition task was used to test their memory effect. Each participant was required to have over 95% accuracy before preceding to matching task. The perceptual learning task was followed, three names were randomly paired with geometric shapes. Participants were required to learn the association and perform a practicing task before they start the formal experimental blocks. They kept practicing until they reached 70% accuracy. Then, they would start the perceptual matching task as in experiment 1a. They finished 6 blocks of perceptual matching trials, each have 120 trials.

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

**Results.** Figure 7 shows  $d$  prime and reaction times of experiment 1c. We conducted same analysis as in Experiment 1a. Our analysis didn't show effect of valence on  $d$  prime,  $F(1.93, 42.56) = 0.23$ ,  $MSE = 0.41$ ,  $p = .791$ ,  $\hat{\eta}_G^2 = .005$ . Neither the effect of valence on RT ( $F(1.63, 35.81) = 0.22$ ,  $MSE = 2,212.71$ ,  $p = .761$ ,  $\hat{\eta}_G^2 = .001$ ) or interaction between valence and matchness on RT ( $F(1.79, 39.43) = 1.20$ ,  $MSE = 1,973.91$ ,  $p = .308$ ,  $\hat{\eta}_G^2 = .005$ ).

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

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

Interesting, we also found the criteria for three conditions also differ, the shapes tagged with good person has the highest criteria (-0.97, [-1.12 -0.82]), followed by shapes

642 tagged with neutral person(-0.96, [-1.09 -0.83]), and then the shapes tagged with bad  
643 person(-1.03, [-1.22 -0.84]). However, pair-wise comparison showed that only showed strong  
644 evidence for the difference between good and bad conditions.

645 ***Reaction time.***

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

654 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et  
655 al., 2020). We estimated separate drift rate ( $v$ ), non-decision time ( $T_0$ ), and boundary  
656 separation ( $a$ ) for each condition. We found that the shapes tagged with good person has  
657 higher drift rate and higher boundary separation than shapes tagged with both neutral and  
658 bad person. Also, the shapes tagged with neutral person has a higher drift rate than  
659 shapes tagged with bad person, but not for the boundary separation. Finally, we found  
660 that shapes tagged with bad person had longer non-decision time (see figure 9)).

661 **Experiment 2: Sequential presenting**

662 Experiment 2 was conducted for two purpose: (1) to further confirm the facilitation  
663 effect of positive moral associations; (2) to test the effect of expectation of occurrence of  
664 each pair. In this experiment, after participant learned the association between labels and  
665 shapes, they were presented a label first and then a shape, they then asked to judge  
666 whether the shape matched the label or not (see (Sui, Sun, Peng, & Humphreys, 2014).

667 Previous studies showed that when the labels presented before the shapes, participants  
668 formed expectations about the shape, and therefore a top-down process were introduced  
669 into the perceptual matching processing. If the facilitation effect of positive moral valence  
670 we found in experiment 1 was mainly drive by top-down processes, this sequential  
671 presenting paradigm may eliminate or attenuate this effect; if, however, the facilitation  
672 effect occurred because of button-up processes, then, similar facilitation effect will appear  
673 even with sequential presenting paradigm.

674 **Method.**

675 ***Participants.***

676 35 participants (17 female, age =  $21.66 \pm 3.03$ ) were recruited. 24 of them had  
677 participated in Experiment 1a (9 male, mean age = 21.9, s.d. = 2.9), and the time gap  
678 between these experiment 1a and experiment 2 is at least six weeks. The results of 1  
679 participants were excluded from analysis because of less than 60% overall accuracy,  
680 remains 34 participants (17 female, age =  $21.74 \pm 3.04$ ).

681 ***Procedure.***

682 In Experiment 2, the sequential presenting makes the matching task much easier than  
683 experiment 1. To avoid ceiling effect on behavioral data, we did a few pilot experiments to  
684 get optimal parameters, i.e., the conditions under which participant have similar accuracy  
685 as in Experiment 1 (around 70 ~ 80% accuracy). In the final procedure, the label (good  
686 person, bad person, or neutral person) was presented for 50 ms and then masked by a  
687 scrambled image for 200 ms. A geometric shape followed the scrambled mask for 50 ms in  
688 a noisy background (which was produced by first decomposing a square with  $\frac{3}{4}$  gray area  
689 and  $\frac{1}{4}$  white area to small squares with a size of  $2 \times 2$  pixels and then re-combine these  
690 small pieces randomly), instead of pure gray background in Experiment 1. After that, a  
691 blank screen was presented 1100 ms, during which participants should press a button to  
692 indicate the label and the shape match the original association or not. Feedback was given,

693 as in study 1. The next trial then started after 700 ~ 1100 ms blank. Other aspects of  
694 study 2 were identical to study 1.

695 ***Data analysis.***

696 Data was analyzed as in study 1a.

697 **Results.**

698 ***NHST.***

699 Figure 10 shows  $d$  prime and reaction times of experiment 2. Less than 0.2% correct  
700 trials with less than 200ms reaction times were excluded.

701 *d prime.*

702 There was evidence for the main effect of valence,  $F(1.83, 60.36) = 14.41$ ,  
703  $MSE = 0.23$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .066$ . Paired t test showed that the Good-Person condition  
704 ( $2.79 \pm 0.17$ ) was with greater  $d$  prime than Netural condition ( $2.21 \pm 0.16$ ,  $t(33) = 4.723$ ,  
705  $p = 0.001$ ) and Bad-person condition ( $2.41 \pm 0.14$ ),  $t(33) = 4.067$ ,  $p = 0.008$ ). There was  
706 no-significant difference between Neutral-person and Bad-person conidition,  $t(33) = -1.802$ ,  
707  $p = 0.185$ .

708 *Reaction time.*

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

711 We found interaction between Matchness and Valence ( $F(1.99, 65.70) = 9.53$ ,  
712  $MSE = 605.36$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .017$ ) and then analyzed the matched trials and  
713 mismatched trials separately, as in experiment 1a. For matched trials, we found the effect  
714 of valence  $F(1.99, 65.76) = 10.57$ ,  $MSE = 1,192.65$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .067$ . Post-hoc  $t$ -tests  
715 revealed that shapes associated with Good Person ( $548 \pm 9.4$ ) were responded faster than  
716 Neutral-Person ( $582 \pm 10.9$ ), ( $t(33) = -3.95$ ,  $p = 0.0011$ ) and Bad Person ( $582 \pm 10.2$ ),  
717  $t(33) = -3.9$ ,  $p = 0.0013$ ). While there was no significant differences between Neutral and

<sup>718</sup> Bad-Person condition ( $t(33) = -0.01, p = 0.999$ ). For non-matched trials, there was no  
<sup>719</sup> significant effect of Valence ( $F(1.99, 65.83) = 0.17, MSE = 489.80, p = .843, \hat{\eta}_G^2 = .001$ ).

<sup>720</sup> **BGLMM.**

<sup>721</sup> *Signal detection theory analysis of accuracy.*

<sup>722</sup> We fitted a Bayesian hierarchical GLM for SDT. The results showed that when the  
<sup>723</sup> shapes were tagged with labels with different moral valence, the sensitivity ( $d'$ ) and criteria  
<sup>724</sup> ( $c$ ) were both influence. For the  $d'$ , we found that the shapes tagged with morally good  
<sup>725</sup> person (2.46, 95% CI[2.21 2.72]) is greater than shapes tagged with moral bad (2.07, 95%  
<sup>726</sup> CI[1.83 2.32]),  $P_{PosteriorComparison} = 1$ . Shape tagged with morally good person is also  
<sup>727</sup> greater than shapes tagged with neutral person (2.23, 95% CI[1.95 2.49]),  
<sup>728</sup>  $P_{PosteriorComparison} = 0.97$ . Also, the shapes tagged with neutral person is greater than  
<sup>729</sup> shapes tagged with morally bad person,  $P_{PosteriorComparison} = 0.92$ .

<sup>730</sup> Interesting, we also found the criteria for three conditions also differ, the shapes  
<sup>731</sup> tagged with good person has the highest criteria (-1.01, [-1.14 -0.88]), followed by shapes  
<sup>732</sup> tagged with neutral person(1.06, [-1.21 -0.92]), and then the shapes tagged with bad  
<sup>733</sup> person(-1.11, [-1.25 -0.97]). However, pair-wise comparison showed that only showed strong  
<sup>734</sup> evidence for the difference between good and bad conditions.

<sup>735</sup> *Reaction times.*

<sup>736</sup> We fitted a Bayesian hierarchical GLM for RTs, with a log-normal distribution as the  
<sup>737</sup> link function. We used the posterior distribution of the regression coefficient to make  
<sup>738</sup> statistical inferences. As in previous studies, the matched conditions are much faster than  
<sup>739</sup> the mismatched trials ( $P_{PosteriorComparison} = .75$ ). We focused on matched trials only, and  
<sup>740</sup> compared different conditions: Good () is not faster than the neutral (),  
<sup>741</sup>  $P_{PosteriorComparison} = .5$ , it was faster than the Bad condition (),  
<sup>742</sup>  $P_{PosteriorComparison} = .88$ . And the neutral condition is faster than the bad condition,  
<sup>743</sup>  $P_{PosteriorComparison} = .95$ . However, the mismatched trials are largely overlapped.

744       **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et  
745 al., 2020). We estimated separate drift rate ( $v$ ), non-decision time ( $T_0$ ), and boundary  
746 separation ( $a$ ) for each condition. We found that the shapes tagged with good person has  
747 higher drift rate and higher boundary separation than shapes tagged with both neutral and  
748 bad person. Also, the shapes tagged with neutral person has a higher drift rate than  
749 shapes tagged with bad person, but not for the boundary separation. Finally, we found  
750 that shapes tagged with bad person had longer non-decision time (see figure  
751 @ref(fig:plot-exp1c -HDDM))).

## 752    Discussion

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

## 759    Experiment 6a: EEG study 1

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

### 762    Method.

#### 763    *Participants.*

764       24 college students (8 female, age =  $22.88 \pm 2.79$ ) participated the current study, all  
765 of them were from Tsinghua University in 2014. Informed consent was obtained from all  
766 participants prior to the experiment according to procedures approved by a local ethics  
767 committee. No participant was excluded from behavioral analysis.

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

772       *Stimuli.*

773       Three geometric shapes (triangle, square and circle, each  $4.6^\circ \times 4.6^\circ$  of visual angle)  
774 were presented at the center of screen for 50 ms after 500ms of fixation ( $0.8^\circ \times 0.8^\circ$  of  
775 visual angle). The association of the three shapes to bad person (" , HuaiRen"), good  
776 person (" , HaoRen") or ordinary person (" , ChangRen") was counterbalanced across  
777 participants. The words bad person, good person or ordinary person ( $3.6^\circ \times 1.6^\circ$ ) was also  
778 displayed at the center fo the screen. Participants had to judge whether the pairings of  
779 label and shape matched (e.g., Does the circle represent a bad person?). The experiment  
780 was run on a PC using E-prime software (version 2.0). These stimuli were displayed on a  
781 22-in CRT monitor ( $1024 \times 768$  at 100Hz). We used backward masking to avoid  
782 over-processing of the moral words, in which a scrambled picture were presented for 900 ms  
783 after the label. Also, to avoid the ceiling effect on accuracy, shapes were presented on a  
784 noisy background based on our pilot studies. The noisy images were made by scrambling a  
785 picture of 3/4gray and 1/4 white at resolution of  $2 \times 2$  pixel.

786       *Procedure.*

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

790       As in experiment 2 (Sui, He, & Humphreys, 2012), subjects first learned the  
791 associations between labels and shapes and then completed a shape-label matching task  
792 (e.g., good person-triangle). In each trial of the matching task, a fixation were first  
793 presented for 500 ms, followed by a 50 ms label; then, a scrambled picture presented 900

794 ms. After the backward mask, the shape were presented on a noisy background for 50ms.  
795 Participant have to response in 1000ms after the presentation of the shape, and finally, a  
796 feedback screen was presented for 500 ms (see figure 1). The inter-trial interval (ITI) were  
797 randomly varied at the range of 1000 ~ 1400 ms.

798 All the stimuli were presented on a gray background (RGB: 127, 127, 127). E-primed  
799 2.0 was used to present stimuli and collect behavioral results. Data were collected and  
800 analyzed when accuracy performance in total reached 60%.

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

## 802 **Results.**

### 803 **NHST.**

804 Only the behavioral results were reported here. Figure 13 shows *d* prime and reaction  
805 times of experiment 6a.

806 *d prime.*

807 We conducted repeated measures ANOVA, with moral valence as independent  
808 variable. The results revealed the main effect of valence ( $F(1.74, 40.05) = 3.76$ ,  
809  $MSE = 0.10$ ,  $p = .037$ ,  $\eta^2_G = .021$ ). Post-hoc analysis revealed that shapes link with Good  
810 person (mean = 3.13, SE = 0.109) is greater than Neutral condition (mean = 2.88, SE =  
811 0.14),  $t = 2.916$ ,  $df = 24$ ,  $p = 0.02$ , p-value adjusted by Tukey method, but the *d* prime  
812 between Good and bad (mean = 3.03, SE = 0.142) ( $t = 1.512$ ,  $df = 24$ ,  $p = 0.3034$ , p-value  
813 adjusted by Tukey method), bad and neutral ( $t = 1.599$ ,  $df = 24$ ,  $p = 0.2655$ , p-value  
814 adjusted by Tukey method) were not significant.

815 *Reaction times.*

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

818 We found intercation between Matchness and Valence ( $F(1.97, 45.20) = 20.45$ ,

819  $MSE = 450.47, p < .001, \hat{\eta}_G^2 = .021$ ) and then analyzed the matched trials and  
 820 mismatched trials separately, as in experiment 2. For matched trials, we found the effect of  
 821 valence  $F(1.97, 45.25) = 32.37, MSE = 522.42, p < .001, \hat{\eta}_G^2 = .078$ . For non-matched  
 822 trials, there was no significant effect of Valence ( $F(1.77, 40.67) = 0.35, MSE = 242.15,$   
 823  $p = .679, \hat{\eta}_G^2 = .000$ ). Post-hoc  $t$ -tests revealed that shapes associated with Good Person  
 824 (mean = 550, SE = 13.8) were responded faster than Neutral-Person (501, SE = 14.7),  
 825 ( $t(24) = -5.171, p = 0.0001$ ) and Bad Person (523, SE = 16.3),  $t(24) = -8.137, p <$   
 826 0.0001),, and Neutral is faster than Bad-Person condition ( $t(32) = -3.282, p = 0.0085$ ).

827 **BGLM.**

828 *Signal detection theory analysis of accuracy.*

829 *Reaction time.*

830 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et  
 831 al., 2020). We estimated separate drift rate ( $v$ ), non-decision time ( $T_0$ ), and boundary  
 832 separation ( $a$ ) for each condition. We found that, similar to experiment 2, the shapes  
 833 tagged with good person has higher drift rate and higher boundary separation than shapes  
 834 tagged with both neutral and bad person, but only for the self-referential condition. Also,  
 835 the shapes tagged with neutral person has a higher drift rate than shapes tagged with bad  
 836 person, but not for the boundary separation, and this effect also exist only for the  
 837 self-referential condition.

838 Interestingly, we found that in both self-referential and other-referential conditions,  
 839 the shapes associated bad valence have higher drift rate and higher boundary separation.  
 840 which might suggest that the shape associated with bad stimuli might be prioritized in the  
 841 non-match trials (see figure 15).

842

## Part 2: interaction between valence and identity

843

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

846

### Experiment 3a

847

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

851

#### Method.

852

##### *Participants.*

853

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

854

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

859

##### *Design.*

860

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

**865      *Stimuli.***

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

**871      *Procedure.***

872      After being fully explained and signed the informed consent, participants were  
873      instructed to chose a name that can represent a stranger with same gender as the  
874      participant themselves, from a common Chinese name pool. Before experiment, the  
875      experimenter explained the meaning of each label to participants. For example, the “good  
876      self” mean the morally good side of themselves, them could imagine the moment when they  
877      do something’s morally applauded, “bad self” means the morally bad side of themselves,  
878      they could also imagine the moment when they doing something morally wrong, and  
879      “neutral self” means the aspect of self that does not related to morality, they could imagine  
880      the moment when they doing something irrelevant to morality. In the same sense, the  
881      “good other”, “bad other”, and “neutral other” means the three different aspects of the  
882      stranger, whose name was chosen before the experiment. Then, the experiment proceeded  
883      as study 1a. Each participant finished 6 blocks, each have 120 trials. The sequence of trials  
884      was pseudo-randomized so that there are 10 matched trials for each condition and 10  
885      non-matched trials for each condition (good self, neutral self, bad self, good other, neutral  
886      other, bad other) for each block.

**887      *Data Analysis.***

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

891       **Results.**

892       **NHST.**

893       Figure 16 shows  $d$  prime and reaction times of experiment 3a. Less than 5% correct  
 894       trials with less than 200ms reaction times were excluded.

895       *d prime.*

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

900       We then conducted separated ANOVA for self-referential and other-referential trials.  
 901       The valence effect was shown for the self-referential conditions,  $F(1.65, 56.25) = 13.98$ ,  
 902        $MSE = 0.31$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .119$ . Post-hoc test revealed that the Good-Self condition  
 903       ( $1.97 \pm 0.14$ ) was with greater  $d$  prime than Neutral condition ( $1.41 \pm 0.12$ ,  $t(34) = 4.505$ ,  
 904        $p = 0.0002$ ), and Bad-self condition ( $1.43 \pm 0.102$ ),  $t(34) = 3.856$ ,  $p = 0.0014$ . There was  
 905       difference between neutral and bad condition,  $t(34) = -0.238$ ,  $p = 0.9694$ . However, no  
 906       effect of valence was found for the other-referential condition  $F(1.98, 67.36) = 0.38$ ,  
 907        $MSE = 0.35$ ,  $p = .681$ ,  $\hat{\eta}_G^2 = .004$ .

908       *Reaction time.*

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

912       For the match trials, we found that the interaction between identity and valence,  
 913        $F(1.72, 58.61) = 3.89$ ,  $MSE = 2,750.19$ ,  $p = .032$ ,  $\hat{\eta}_G^2 = .019$ , as well as the main effect of  
 914       valence  $F(1.98, 67.34) = 35.76$ ,  $MSE = 1,127.25$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .079$ , but not the effect of  
 915       identity  $F(1, 34) = 0.20$ ,  $MSE = 3,507.14$ ,  $p = .660$ ,  $\hat{\eta}_G^2 = .001$ . As for the  $d$  prime, we

916 separated analyzed the self-referential and other-referential trials. For the Self-referential  
 917 trials, we found the main effect of valence,  $F(1.80, 61.09) = 30.39$ ,  $MSE = 1,584.53$ ,  
 918  $p < .001$ ,  $\hat{\eta}_G^2 = .159$ ; for the other-referential trials, the effect of valence is weaker,  
 919  $F(1.86, 63.08) = 2.85$ ,  $MSE = 2,224.30$ ,  $p = .069$ ,  $\hat{\eta}_G^2 = .024$ . We then focused on the self  
 920 conditions: the good-self condition ( $713 \pm 12$ ) is faster than neutral- ( $776 \pm 11.8$ ),  $t(34) =$   
 921  $-7.396$ ,  $p < .0001$ , and bad-self ( $772 \pm 10.1$ ) conditions,  $t(34) = -5.66$ ,  $p < .0001$ . But  
 922 there is not difference between neutral- and bad-self conditions,  $t(34) = 0.481$ ,  $p = 0.881$ .

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

## 927 **BGLM.**

928 *Signal detection theory analysis of accuracy.*

929 We found that the  $d$  prime is greater when shapes were associated with good self  
 930 condition than with neutral self or bad self, but shapes associated with bad self and neutral  
 931 self didn't show differences. Comparing the self vs other under three condition revealed  
 932 that shapes associated with good self is greater than with good other, but with a weak  
 933 evidence. In contrast, for both neutral and bad valence condition, shapes associated with  
 934 other had greater  $d$  prime than with self.

935 *Reaction time.*

936 In reaction times, we found that same trends in the match trials as in the RT: while  
 937 the shapes associated with good self was greater than with good other (log mean diff =  
 938  $-0.02858$ , 95%HPD[-0.070898, 0.0154]), the direction is reversed for neutral and negative  
 939 condition. see Figure 17

940 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et  
 941 al., 2020). We estimated separate drift rate ( $v$ ), non-decision time ( $T_0$ ), and boundary

942 separation (*a*) for each condition. We found that the shapes tagged with good person has  
943 higher drift rate and higher boundary separation than shapes tagged with both neutral and  
944 bad person. Also, the shapes tagged with neutral person has a higher drift rate than  
945 shapes tagged with bad person, but not for the boundary separation. Finally, we found  
946 that shapes tagged with bad person had longer non-decision time (see figure 18)).

947 **Experiment 3b**

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

953 **Participants.**

954 Study 3b were finished in 2017, at that time we have calculated that the effect size  
955 (Cohen's *d*) of good-person (or good-self) vs. bad-person (or bad-other) was between  $0.47 \sim 0.53$ , based on the data from Tsinghua community in study 1a, 1b, 2, 3a, 4a, and 4b. Based  
956 on this effect size, we estimated that 54 participants would allow we to detect the effect  
957 size of Cohen's  $= 0.5$  with 95% power and alpha = 0.05, using G\*power 3.192 (Faul,  
958 Erdfelder, Buchner, & Lang, 2009). Therefore, we planned to stop after we arrived this  
959 number. During the data collected at Wenzhou University, 61 participants (45 females; 19  
960 to 25 years of age, age =  $20.42 \pm 1.77$ ) came to the testing room and we tested all of them  
961 during a single day. All participants were right-handed, and all had normal or  
962 corrected-to-normal vision. Informed consent was obtained from all participants prior to  
963 the experiment according to procedures approved by a local ethics committee. 4  
964 participants' data were excluded from analysis because their over all accuracy was lower  
965 than 60%, 1 more participant was excluded because of zero hit rate for one condition,  
966 leaving 56 participants (43 females; 19 to 25 years old, age =  $20.27 \pm 1.60$ ).

***Design.***

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

***Stimuli.***

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

***Procedure.***

In this experiment, participants finished two matching tasks, i.e., self-matching task, and other-matching task. In the self-matching task, participants first associate the three aspects of self to three different shapes, and then perform the matching task. In the other-matching task, participants first associate the three aspects of the stranger to three different shapes, and then perform the matching task. The order of self-task and other-task are counter-balanced among participants. Different from experiment 3a, after presenting the stimuli pair for 100ms, participant has 1900 ms to response, and they feedback with both accuracy and reaction time. 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.

994        ***Data Analysis.***

995        Same as experiment 3a.

996        **Results.**

997        ***NHST.***

998        Figure 19 shows  $d$  prime and reaction times of experiment 3b. Less than 5% correct  
 999        trials with less than 200ms reaction times were excluded.

1000        *d prime.*

1001        There was no evidence for the main effect of valence,  $F(1.92, 105.43) = 1.90$ ,

1002         $MSE = 0.33$ ,  $p = .157$ ,  $\hat{\eta}_G^2 = .005$ , but we found a main effect of self-relevance,

1003         $F(1, 55) = 4.65$ ,  $MSE = 0.89$ ,  $p = .035$ ,  $\hat{\eta}_G^2 = .017$ , as well as the interaction,

1004         $F(1.90, 104.36) = 5.58$ ,  $MSE = 0.26$ ,  $p = .006$ ,  $\hat{\eta}_G^2 = .011$ .

1005        We then conducted separated ANOVA for self-referential and other-referential trials.

1006        The valence effect was shown for the self-referential conditions,  $F(1.75, 96.42) = 6.73$ ,

1007         $MSE = 0.30$ ,  $p = .003$ ,  $\hat{\eta}_G^2 = .037$ . Post-hoc test revealed that the Good-Self condition

1008         $(2.15 \pm 0.12)$  was with greater  $d$  prime than Neutral condition  $(1.83 \pm 0.12$ ,  $t(34) = 3.36$ ,

1009         $p = 0.0031$ ), and Bad-self condition  $(1.87 \pm 0.12)$ ,  $t(34) = 2.955$ ,  $p = 0.01$ . There was

1010        difference between neutral and bad condition,  $t(34) = -0.039$ ,  $p = 0.914$ . However, no effect

1011        of valence was found for the other-referential condition  $F(1.93, 105.97) = 0.61$ ,

1012         $MSE = 0.31$ ,  $p = .539$ ,  $\hat{\eta}_G^2 = .002$ .

1013        *Reaction time.*

1014        We found interaction between Matchness and Valence ( $F(1.86, 102.47) = 15.44$ ,

1015         $MSE = 3, 112.78$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .006$ ) and then analyzed the matched trials and

1016        nonmatch trials separately, as in previous experiments.

1017        For the match trials, we found that the interaction between identity and valence,

1018         $F(1.67, 92.11) = 6.14$ ,  $MSE = 6, 472.48$ ,  $p = .005$ ,  $\hat{\eta}_G^2 = .009$ , as well as the main effect of

1019 valence  $F(1.88, 103.65) = 24.25$ ,  $MSE = 5,994.25$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .038$ , but not the effect  
 1020 of identity  $F(1, 55) = 48.49$ ,  $MSE = 25,892.59$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .153$ . As for the  $d$  prime,  
 1021 we separated analyzed the self-referential and other-referential trials. For the  
 1022 Self-referential trials, we found the main effect of valence,  $F(1.66, 91.38) = 23.98$ ,  
 1023  $MSE = 6,965.61$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .100$ ; for the other-referential trials, the effect of valence  
 1024 is weaker,  $F(1.89, 103.94) = 5.96$ ,  $MSE = 5,589.90$ ,  $p = .004$ ,  $\hat{\eta}_G^2 = .014$ . We then focused  
 1025 on the self conditions: the good-self condition ( $713 \pm 12$ ) is faster than neutral- ( $776 \pm$   
 1026  $11.8$ ),  $t(34) = -7.396$ ,  $p < .0001$ , and bad-self ( $772 \pm 10.1$ ) conditions,  $t(34) = -5.66$ ,  $p <$   
 1027  $.0001$ . But there is not difference between neutral- and bad-self conditions,  $t(34) = 0.481$ ,  $p$   
 1028 = 0.881.

1029 For the nonmatch trials, we didn't found any strong effect: identity,  $F(1, 55) = 10.31$ ,  
 1030  $MSE = 24,590.52$ ,  $p = .002$ ,  $\hat{\eta}_G^2 = .035$ , valence  $F(1.98, 108.63) = 20.57$ ,  $MSE = 2,847.51$ ,  
 1031  $p < .001$ ,  $\hat{\eta}_G^2 = .016$ , or interaction between the two  $F(1.93, 106.25) = 35.51$ ,  
 1032  $MSE = 1,939.88$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .019$ .

### 1033 **BGLM.**

1034 *Signal detection theory analysis of accuracy.*

1035 We found that the  $d$  prime is greater when shapes were associated with good self  
 1036 condition than with neutral self or bad self, but shapes associated with bad self and neutral  
 1037 self didn't show differences. comparing the self vs other under three condition revealed that  
 1038 shapes associated with good self is greater than with good other, but with a weak evidence.  
 1039 In contrast, for both neutral and bad valence condition, shapes associated with other had  
 1040 greater  $d$  prime than with self.

1041 *Reaction time.*

1042 In reaction times, we found that same trends in the match trials as in the RT: while  
 1043 the shapes associated with good self was greater than with good other (log mean diff =  
 1044 -0.02858, 95%HPD[-0.070898, 0.0154]), the direction is reversed for neutral and negative

1045 condition. see Figure 20

1046 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et  
1047 al., 2020). We estimated separate drift rate ( $v$ ), non-decision time ( $T_0$ ), and boundary  
1048 separation ( $a$ ) for each condition. We found that, similar to experiment 3a, the shapes  
1049 tagged with good person has higher drift rate and higher boundary separation than shapes  
1050 tagged with both neutral and bad person, but only for the self-referential condition. Also,  
1051 the shapes tagged with neutral person has a higher drift rate than shapes tagged with bad  
1052 person, but not for the boundary separation, and this effect also exist only for the  
1053 self-referential condition.

1054 Interestingly, we found that in both self-referential and other-referential conditions,  
1055 the shapes associated bad valence have higher drift rate and higher boundary separation.  
1056 which might suggest that the shape associated with bad stimuli might be prioritized in the  
1057 non-match trials (see figure 21)).

## 1058 **Experiment 6b**

1059 Experiment 6b was conducted to study the neural correlates of the prioritization  
1060 effect of positive self, i.e., the neural underlying of the behavioral effect found int  
1061 experiment 3a. However, as in experiment 6a, the procedure of this experiment was  
1062 modified to adopted to ERP experiment.

### 1063 **Method.**

#### 1064 ***Participants.***

1065 23 college students (8 female, age =  $22.86 \pm 2.47$ ) participated the current study, all  
1066 of them were recruited from Tsinghua University in 2016. Informed consent was obtained  
1067 from all participants prior to the experiment according to procedures approved by a local  
1068 ethics committee. For day 1's data, 1 participant was excluded from the current analysis  
1069 because of lower than 60% overall accuracy, remaining 22 participants (8 female, age =

1070 22.76 ± 2.49). For day 2's data, one participant dropped out, leaving 22 participants (9  
1071 female, age = 23.05 ± 2.46), all of them has overall accuracy higher than 60%.

1072 ***Design.***

1073 The experimental design of this experiment is same as experiment 3: a 2 × 3 × 2  
1074 within-subject design with self-relevance (self-relevant vs. other-relevant), moral valence  
1075 (good, neutral, and bad) and matchness between shape and label (match vs. mismatch) as  
1076 within-subject variables.

1077 ***Stimuli.***

1078 As in experiment 3a, 6 shapes were included (triangle, square, circle, trapezoid,  
1079 diamond, regular pentagon), as well as 6 labels (good self, neutral self, bad self, good  
1080 person, bad person, neutral person). To match the concreteness of the label, we asked  
1081 participant to chosen an unfamiliar name of their own gender to be the stranger.

1082 ***Procedure.***

1083 The procedure was similar to Experiment 2 and 6a. Subjects first learned the  
1084 associations between labels and shapes and then completed a shape-label matching task. In  
1085 each trial of the matching task, a fixation were first presented for 500 ms, followed by a 50  
1086 ms label; then, a scrambled picture presented 900 ms. After the backward mask, the shape  
1087 were presented on a noisy background for 50ms. Participant have to response in 1000ms  
1088 after the presentation of the shape, and finally, a feedback screen was presented for 500 ms.  
1089 The inter-trial interval (ITI) were randomly varied at the range of 1000 ~ 1400 ms.

1090 All the stimuli were presented on a gray background (RGB: 127, 127, 127). E-primed  
1091 2.0 was used to present stimuli and collect behavioral results. Data were collected and  
1092 analyzed when accuracy performance in total reached 60%.

1093 Because learning 6 associations was more difficult than 3 associations and participant  
1094 might have low accuracy (see experiment 3a), the current study had extended to a two-day

1095 paradigm to maximizing the accurate trials that can be used in EEG data. At the first day,  
1096 participants learnt the associations and finished 9 blocks of the matching task, each had  
1097 120 trials, without EEG recording. That is, each condition has 90 trials.

1098 Participants came back to lab at the second day and finish the same task again, with  
1099 EEG recorded. Before the EEG experiment, each participant finished a practice session  
1100 again, if their accuracy is equal or higher than 85%, they start the experiment (one  
1101 participant used lower threshold 75%). Each participant finished 18 blocks, each has 90  
1102 trials. One participant finished additional 6 blocks because of high error rate at the  
1103 beginning, another two participant finished addition 3 blocks because of the technique  
1104 failure in recording the EEG data. To increase the number of trials that can be used for  
1105 EEG data analysis, matched trials has twice number as mismatched trials, therefore, for  
1106 matched trials each participants finished 180 trials for each condition, for mismatched  
1107 trials, each conditions has 90 trials.

1108 ***Data Analysis.***

1109 Same as experiment 3a.

1110 **Results of Day 1.**

1111 ***NHST.***

1112 Figure 22 shows  $d$  prime and reaction times of experiment 3b. Less than 5% correct  
1113 trials with less than 200ms reaction times were excluded.

1114 ***d prime.***

1115 There was no evidence for the main effect of valence,  $F(1.91, 40.20) = 11.98$ ,  
1116  $MSE = 0.15$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .040$ , but we found a main effect of self-relevance,  
1117  $F(1, 21) = 1.21$ ,  $MSE = 0.20$ ,  $p = .284$ ,  $\hat{\eta}_G^2 = .003$ , as well as the interaction,  
1118  $F(1.28, 26.90) = 12.88$ ,  $MSE = 0.21$ ,  $p = .001$ ,  $\hat{\eta}_G^2 = .041$ .

1119 We then conducted separated ANOVA for self-referential and other-referential trials.

1120 The valence effect was shown for the self-referential conditions,  $F(1.73, 36.42) = 29.31$ ,  
 1121  $MSE = 0.14$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .147$ . Post-hoc test revealed that the Good-Self condition  
 1122 ( $2.15 \pm 0.12$ ) was with greater  $d$  prime than Neutral condition ( $1.83 \pm 0.12$ ,  $t(34) = 3.36$ ,  
 1123  $p = 0.0031$ ), and Bad-self condition ( $1.87 \pm 0.12$ ),  $t(34) = 2.955$ ,  $p = 0.01$ . There was  
 1124 difference between neutral and bad condition,  $t(34) = -0.039$ ,  $p = 0.914$ . However, no effect  
 1125 of valence was found for the other-referential condition  $F(1.75, 36.72) = 0.00$ ,  $MSE = 0.18$ ,  
 1126  $p = .999$ ,  $\hat{\eta}_G^2 = .000$ .

1127 *Reaction time.*

1128 We found interaction between Matchness and Valence ( $F(1.79, 37.63) = 4.07$ ,  
 1129  $MSE = 704.90$ ,  $p = .029$ ,  $\hat{\eta}_G^2 = .003$ ) and then analyzed the matched trials and nonmatch  
 1130 trials separately, as in previous experiments.

1131 For the match trials, we found that the interaction between identity and valence,  
 1132  $F(1.72, 36.16) = 4.55$ ,  $MSE = 1,560.90$ ,  $p = .022$ ,  $\hat{\eta}_G^2 = .015$ , as well as the main effect of  
 1133 valence  $F(1.93, 40.55) = 9.83$ ,  $MSE = 1,951.84$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .044$ , but not the effect of  
 1134 identity  $F(1, 21) = 4.87$ ,  $MSE = 2,032.05$ ,  $p = .039$ ,  $\hat{\eta}_G^2 = .012$ . As for the  $d$  prime, we  
 1135 separated analyzed the self-referential and other-referential trials. For the Self-referential  
 1136 trials, we found the main effect of valence,  $F(1.92, 40.38) = 14.48$ ,  $MSE = 1,647.20$ ,  
 1137  $p < .001$ ,  $\hat{\eta}_G^2 = .112$ ; for the other-referential trials, the effect of valence is weaker,  
 1138  $F(1.79, 37.50) = 1.04$ ,  $MSE = 1,842.07$ ,  $p = .356$ ,  $\hat{\eta}_G^2 = .008$ . We then focused on the self  
 1139 conditions: the good-self condition ( $713 \pm 12$ ) is faster than neutral- ( $776 \pm 11.8$ ),  $t(34) =$   
 1140  $-7.396$ ,  $p < .0001$ , and bad-self ( $772 \pm 10.1$ ) conditions,  $t(34) = -5.66$ ,  $p < .0001$ . But  
 1141 there is not difference between neutral- and bad-self conditions,  $t(34) = 0.481$ ,  $p = 0.881$ .

1142 For the nonmatch trials, we didn't found any strong effect: identity,  $F(1, 21) = 2.76$ ,  
 1143  $MSE = 1,718.93$ ,  $p = .112$ ,  $\hat{\eta}_G^2 = .006$ , valence  $F(1.61, 33.77) = 3.81$ ,  $MSE = 1,532.21$ ,  
 1144  $p = .041$ ,  $\hat{\eta}_G^2 = .012$ , or interaction between the two  $F(1.90, 39.97) = 2.23$ ,  $MSE = 720.80$ ,  
 1145  $p = .123$ ,  $\hat{\eta}_G^2 = .004$ .

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

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

*Reaction time.*

In reaction times, we found that same trends in the match trials as in the RT: while the shapes associated with good self was greater than with good other (log mean diff = -0.02858, 95%HPD[-0.070898, 0.0154]), the direction is reversed for neutral and negative condition. see Figure 23

**HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et al., 2020). We estimated separate drift rate ( $v$ ), non-decision time ( $T_0$ ), and boundary separation ( $a$ ) for each condition. We found that, similar to experiment 3a, the shapes tagged with good person has higher drift rate and higher boundary separation than shapes tagged with both neutral and bad person, but only for the self-referential condition. Also, the shapes tagged with neutral person has a higher drift rate than shapes tagged with bad person, but not for the boundary separation, and this effect also exist only for the self-referential condition.

Interestingly, we found that in both self-referential and other-referential conditions, the shapes associated bad valence have higher drift rate and higher boundary separation. which might suggest that the shape associated with bad stimuli might be prioritized in the non-match trials (see figure 24).

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

In this part, we reported two studies in which the moral valence or the self-referential processing is not task-relevant. We are interested in testing whether the task-relevance will eliminate the effect observed in previous experiment.

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

In part two (experiment 3a and 3b), participants learned the association between self and moral valence directly. In Experiment 4a, we examined whether the interaction between moral valence and identity occur even when one of the variable was irrelevant to the task. In experiment 4a, participants learnt associations between shapes and self/other labels, then made perceptual match judgments only about the self or other conditions labels and shapes (cf. Sui et al. (2012)). However, we presented labels of different moral valence in the shapes, which means that the moral valence factor become task irrelevant. If the binding between moral good and self is intrinsic and automatic, then we will observe that facilitating effect of moral good for self conditions, but not for other conditions.

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

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

***Design.***

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

<sub>1221</sub> 120 trials for 2017 data. Due to procedure error, the data collected in 2015 in Tsinghua  
<sub>1222</sub> community only have 60 trials for each block, i.e., 30 trials per condition.

<sub>1223</sub> As in study 3a, before each task, the instruction showed the meaning of each label to  
<sub>1224</sub> participants. The self-matching task and other-matching task were randomized between  
<sub>1225</sub> participants. Each participant finished 6 blocks, each have 120 trials.

<sub>1226</sub> ***Data Analysis.***

<sub>1227</sub> Same as experiment 3a.

<sub>1228</sub> **Results.**

<sub>1229</sub> ***NHST.***

<sub>1230</sub> Figure 25 shows  $d$  prime and reaction times of experiment 3a. Less than 5% correct  
<sub>1231</sub> trials with less than 200ms reaction times were excluded.

<sub>1232</sub>  $d$  prime.

<sub>1233</sub> There was no evidence for the main effect of valence,  $F(1.93, 111.66) = 0.53$ ,  
<sub>1234</sub>  $MSE = 0.12$ ,  $p = .581$ ,  $\hat{\eta}_G^2 = .000$ , but we found a main effect of self-relevance,  
<sub>1235</sub>  $F(1, 58) = 121.04$ ,  $MSE = 0.48$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .189$ , as well as the interaction,  
<sub>1236</sub>  $F(1.99, 115.20) = 4.12$ ,  $MSE = 0.14$ ,  $p = .019$ ,  $\hat{\eta}_G^2 = .004$ .

<sub>1237</sub> We then conducted separated ANOVA for self-referential and other-referential trials.

<sub>1238</sub> The valence effect was shown for the self-referential conditions,  $F(1.95, 112.92) = 3.01$ ,  
<sub>1239</sub>  $MSE = 0.15$ ,  $p = .055$ ,  $\hat{\eta}_G^2 = .008$ . Post-hoc test revealed that the Good-Self condition  
<sub>1240</sub> ( $2.15 \pm 0.12$ ) was with greater  $d$  prime than Neutral condition ( $1.83 \pm 0.12$ ,  $t(34) = 3.36$ ,  
<sub>1241</sub>  $p = 0.0031$ ), and Bad-self condition ( $1.87 \pm 0.12$ ),  $t(34) = 2.955$ ,  $p = 0.01$ . There was  
<sub>1242</sub> difference between neutral and bad condition,  $t(34) = -0.039$ ,  $p = 0.914$ . However, no effect  
<sub>1243</sub> of valence was found for the other-referential condition  $F(1.98, 114.61) = 1.75$ ,  
<sub>1244</sub>  $MSE = 0.10$ ,  $p = .179$ ,  $\hat{\eta}_G^2 = .003$ .

<sub>1245</sub> Reaction time.

1246 We found interaction between Matchness and Valence ( $F(1.94, 112.64) = 0.84$ ,  
 1247  $MSE = 465.35, p = .432, \hat{\eta}_G^2 = .000$ ) and then analyzed the matched trials and nonmatch  
 1248 trials separately, as in previous experiments.

1249 For the match trials, we found that the interaction between identity and valence,  
 1250  $F(1.90, 110.18) = 4.41, MSE = 465.91, p = .016, \hat{\eta}_G^2 = .003$ , as well as the main effect of  
 1251 valence  $F(1.98, 114.82) = 0.94, MSE = 606.30, p = .392, \hat{\eta}_G^2 = .001$ , but not the effect of  
 1252 identity  $F(1, 58) = 124.15, MSE = 4,037.53, p < .001, \hat{\eta}_G^2 = .257$ . As for the  $d$  prime, we  
 1253 separated analyzed the self-referential and other-referential trials. For the Self-referential  
 1254 trials, we found the main effect of valence,  $F(1.97, 114.32) = 6.29, MSE = 367.25$ ,  
 1255  $p = .003, \hat{\eta}_G^2 = .006$ ; for the other-referential trials, the effect of valence is weaker,  
 1256  $F(1.95, 112.89) = 0.35, MSE = 699.50, p = .699, \hat{\eta}_G^2 = .001$ . We then focused on the self  
 1257 conditions: the good-self condition ( $713 \pm 12$ ) is faster than neutral- ( $776 \pm 11.8$ ),  $t(34) =$   
 1258  $-7.396, p < .0001$ , and bad-self ( $772 \pm 10.1$ ) conditions,  $t(34) = -5.66, p < .0001$ . But  
 1259 there is not difference between neutral- and bad-self conditions,  $t(34) = 0.481, p = 0.881$ .

1260 For the nonmatch trials, we didn't found any strong effect: identity,  $F(1, 58) = 0.16$ ,  
 1261  $MSE = 1,547.37, p = .692, \hat{\eta}_G^2 = .000$ , valence  $F(1.96, 113.52) = 0.68, MSE = 390.26$ ,  
 1262  $p = .508, \hat{\eta}_G^2 = .000$ , or interaction between the two  $F(1.90, 110.27) = 0.04$ ,  
 1263  $MSE = 585.80, p = .953, \hat{\eta}_G^2 = .000$ .

1264 **BGLM.**

1265 *Signal detection theory analysis of accuracy.*

1266 We found that the  $d$  prime is greater when shapes were associated with good self  
 1267 condition than with neutral self or bad self, but shapes associated with bad self and neutral  
 1268 self didn't show differences. comparing the self vs other under three condition revealed that  
 1269 shapes associated with good self is greater than with good other, but with a weak evidence.  
 1270 In contrast, for both neutral and bad valence condition, shapes associated with other had  
 1271 greater  $d$  prime than with self.

1272       *Reaction time.*

1273       In reaction times, we found that same trends in the match trials as in the RT: while  
1274       the shapes associated with good self was greater than with good other (log mean diff =  
1275       -0.02858, 95%HPD[-0.070898, 0.0154]), the direction is reversed for neutral and negative  
1276       condition. see Figure 26

1277       **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et  
1278       al., 2020). We estimated separate drift rate ( $v$ ), non-decision time ( $T_0$ ), and boundary  
1279       separation ( $a$ ) for each condition. We found that the shapes tagged with good person has  
1280       higher drift rate and higher boundary separation than shapes tagged with both neutral and  
1281       bad person. Also, the shapes tagged with neutral person has a higher drift rate than  
1282       shapes tagged with bad person, but not for the boundary separation. Finally, we found  
1283       that shapes tagged with bad person had longer non-decision time (see figure 27)).

1284       **Experiment 4b: Morality as task-irrelevant variable**

1285       In study 4b, we changed the role of valence and identity in task. In this experiment,  
1286       participants learn the association between moral valence and the made perceptual match  
1287       judgments to associations between different moral valence and shapes as in study 1-3.  
1288       Different from experiment 1 ~ 3, we made put the labels of “self/other” in the shapes so  
1289       that identity served as an task irrelevant variable. As in experiment 4b, we also  
1290       hypothesized that the intrinsic binding between morally good and self will enhance the  
1291       performance of good self condition, even identity is irrelevant to the task.

1292       **Method.**

1293       **Participants.**

1294       53 participants (39 female, age =  $20.57 \pm 1.81$ ) participated the current study, 34 of  
1295       them were from Tsinghua University in 2015, 19 were from Wenzhou University  
1296       participated in 2017. All participants were right-handed, and all had normal or

1297 corrected-to-normal vision. Informed consent was obtained from all participants prior to  
1298 the experiment according to procedures approved by a local ethics committee. The data  
1299 from 8 participants from Wenzhou site were excluded from analysis because their accuracy  
1300 was close to chance (< 0.6). The results for the remaining 45 participants (33 female, age  
1301 = 20.78 ± 1.76) were analyzed and reported.

1302 ***Design.***

1303 As in Experiment 3, a 2×3×2 within-subject design was used. The first variable was  
1304 self-relevance (self and stranger associations); the second variable was moral valence (good,  
1305 neutral and bad associations); the third variable was the matching between shape and label  
1306 (matching vs. non-match for the personal association). However, in this the task,  
1307 participants only learn the association between two geometric shapes and two labels (self  
1308 and other), i.e., only self-relevance were related to the task. The moral valence  
1309 manipulation was achieved by embedding the personal label of the labels in the geometric  
1310 shapes, see below. For simplicity, the trials where shapes where paired with self and with a  
1311 word of “good person” inside were shorted as good-self condition, similarly, the trials where  
1312 shapes paired with the self and with a word of “bad person” inside were shorted as bad-self  
1313 condition. Hence, we also have six conditions: good-self, neutral-self, bad-self, good-other,  
1314 neutral-other, and bad-other.

1315 ***Stimuli.***

1316 2 shapes were included (circle, square) and each appeared above a central fixation  
1317 cross with the personal label appearing below. However, the shapes were not empty but  
1318 with a two-Chinese-character word in the middle, the word was one of three labels with  
1319 different moral valence: “good person”, “bad person” and “neutral person”. Before the  
1320 experiment, participants learned the self/other association, and were informed to only  
1321 response to the association between shapes’ configure and the labels below the fixation, but  
1322 ignore the words within shapes. Besides the behavioral experiments, participants from

1323 Tsinghua community also finished questionnaires as Experiments 3, and participants from  
1324 Wenzhou community finished a series of questionnaire as the other experiment finished in  
1325 Wenzhou.

1326 ***Procedure.***

1327 The procedure was similar to Experiment 1. There were 6 blocks of trial, each with  
1328 120 trials for 2017 data. Due to procedure error, the data collected in 2015 in Tsinghua  
1329 community only have 60 trials for each block, i.e., 30 trials per condition.

1330 As in study 3a, before each task, the instruction showed the meaning of each label to  
1331 participants. The self-matching task and other-matching task were randomized between  
1332 participants. Each participant finished 6 blocks, each have 120 trials.

1333 ***Data Analysis.***

1334 Same as experiment 3a.

1335 **Results.**

1336 ***NHST.***

1337 Figure 28 shows  $d$  prime and reaction times of experiment 3a. Less than 5% correct  
1338 trials with less than 200ms reaction times were excluded.

1339  $d$  prime.

1340 There was no evidence for the main effect of valence,  $F(1.59, 69.94) = 2.34$ ,  
1341  $MSE = 0.48$ ,  $p = .115$ ,  $\hat{\eta}_G^2 = .010$ , but we found a main effect of self-relevance,  
1342  $F(1, 44) = 0.00$ ,  $MSE = 0.08$ ,  $p = .994$ ,  $\hat{\eta}_G^2 = .000$ , as well as the interaction,  
1343  $F(1.96, 86.41) = 0.53$ ,  $MSE = 0.10$ ,  $p = .585$ ,  $\hat{\eta}_G^2 = .001$ .

1344 We then conducted separated ANOVA for self-referential and other-referential trials.  
1345 The valence effect was shown for the self-referential conditions,  $F(1.75, 76.86) = 3.08$ ,  
1346  $MSE = 0.25$ ,  $p = .058$ ,  $\hat{\eta}_G^2 = .017$ . Post-hoc test revealed that the Good-Self condition  
1347 ( $2.15 \pm 0.12$ ) was with greater  $d$  prime than Neutral condition ( $1.83 \pm 0.12$ ,  $t(34) = 3.36$ ,

<sup>1348</sup>  $p = 0.0031$ ), and Bad-self condition ( $1.87 \pm 0.12$ ),  $t(34) = 2.955$ ,  $p = 0.01$ . There was  
<sup>1349</sup> difference between neutral and bad condition,  $t(34) = -0.039$ ,  $p = 0.914$ . However, no effect  
<sup>1350</sup> of valence was found for the other-referential condition  $F(1.63, 71.50) = 1.07$ ,  $MSE = 0.33$ ,  
<sup>1351</sup>  $p = .336$ ,  $\hat{\eta}_G^2 = .006$ .

<sup>1352</sup> *Reaction time.*

<sup>1353</sup> We found interaction between Matchness and Valence ( $F(1.87, 82.50) = 18.58$ ,  
<sup>1354</sup>  $MSE = 1,291.12$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .023$ ) and then analyzed the matched trials and  
<sup>1355</sup> nonmatch trials separately, as in previous experiments.

<sup>1356</sup> For the match trials, we found that the interaction between identity and valence,  
<sup>1357</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>1358</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>1359</sup> identity  $F(1, 44) = 0.23$ ,  $MSE = 263.26$ ,  $p = .632$ ,  $\hat{\eta}_G^2 = .000$ . As for the  $d$  prime, we  
<sup>1360</sup> separated analyzed the self-referential and other-referential trials. For the Self-referential  
<sup>1361</sup> trials, we found the main effect of valence,  $F(1.74, 76.48) = 13.69$ ,  $MSE = 1,732.08$ ,  
<sup>1362</sup>  $p < .001$ ,  $\hat{\eta}_G^2 = .079$ ; for the other-referential trials, the effect of valence is weaker,  
<sup>1363</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>1364</sup> conditions: the good-self condition ( $713 \pm 12$ ) is faster than neutral- ( $776 \pm 11.8$ ),  $t(34) =$   
<sup>1365</sup>  $-7.396$ ,  $p < .0001$ , and bad-self ( $772 \pm 10.1$ ) conditions,  $t(34) = -5.66$ ,  $p < .0001$ . But  
<sup>1366</sup> there is not difference between neutral- and bad-self conditions,  $t(34) = 0.481$ ,  $p = 0.881$ .

<sup>1367</sup> For the nonmatch trials, we didn't found any strong effect: identity,  $F(1, 44) = 1.96$ ,  
<sup>1368</sup>  $MSE = 319.47$ ,  $p = .169$ ,  $\hat{\eta}_G^2 = .001$ , valence  $F(1.69, 74.54) = 6.59$ ,  $MSE = 886.19$ ,  
<sup>1369</sup>  $p = .004$ ,  $\hat{\eta}_G^2 = .010$ , or interaction between the two  $F(1.88, 82.57) = 0.31$ ,  $MSE = 316.96$ ,  
<sup>1370</sup>  $p = .718$ ,  $\hat{\eta}_G^2 = .000$ .

<sup>1371</sup> **BGLM.**

<sup>1372</sup> *Signal detection theory analysis of accuracy.*

<sup>1373</sup> We found that the  $d$  prime is greater when shapes were associated with good self

1374 condition than with neutral self or bad self, but shapes associated with bad self and neutral  
1375 self didn't show differences. comparing the self vs other under three condition revealed that  
1376 shapes associated with good self is greater than with good other, but with a weak evidence.  
1377 In contrast, for both neutral and bad valence condition, shapes associated with other had  
1378 greater  $d$  prime than with self.

1379 *Reaction time.*

1380 In reaction times, we found that same trends in the match trials as in the RT: while  
1381 the shapes associated with good self was greater than with good other (log mean diff =  
1382 -0.02858, 95%HPD[-0.070898, 0.0154]), the direction is reversed for neutral and negative  
1383 condition. see Figure 29

1384 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et  
1385 al., 2020). We estimated separate drift rate ( $v$ ), non-decision time ( $T_0$ ), and boundary  
1386 separation ( $a$ ) for each condition. We found that the shapes tagged with good person has  
1387 higher drift rate and higher boundary separation than shapes tagged with both neutral and  
1388 bad person. Also, the shapes tagged with neutral person has a higher drift rate than  
1389 shapes tagged with bad person, but not for the boundary separation. Finally, we found  
1390 that shapes tagged with bad person had longer non-decision time (see figure 30)).

1391

## Results

1392 **Effect of moral valence**

1393 In this part, we synthesized results from experiment 1a, 1b, 1c, 2, 5 and 6a. Data  
1394 from 192 participants were included in these analyses. We found differences between  
1395 positive and negative conditions on RT was Cohen's  $d = -0.58 \pm 0.06$ , 95% CI [-0.70 -0.47];  
1396 on  $d'$  was Cohen's  $d = 0.24 \pm 0.05$ , 95% CI [0.15 0.34]. The effect was also observed  
1397 between positive and neutral condition, RT: Cohen's  $d = -0.44 \pm 0.10$ , 95% CI [-0.63  
1398 -0.25];  $d'$ : Cohen's  $d = 0.31 \pm 0.07$ , 95% CI [0.16 0.45]. And the difference between neutral

<sup>1399</sup> and bad conditions are not significant, RT: Cohen's  $d = 0.15 \pm 0.07$ , 95% CI [0.00 0.30];  
<sup>1400</sup>  $d'$ : Cohen's  $d = 0.07 \pm 0.07$ , 95% CI [-0.08 0.21]. See Figure 31 left panel.

<sup>1401</sup> **Interaction between valence and self-reference**

<sup>1402</sup> In this part, we combined the experiments that explicitly manipulated the  
<sup>1403</sup> self-reference and valence, which includes 3a, 3b, 6b, 7a, and 7b. For the positive versus  
<sup>1404</sup> negative contrast, data were from five experiments with 178 participants; for positive  
<sup>1405</sup> versus neutral and neutral versus negative contrasts, data were from three experiments ( (

<sup>1406</sup> 3a, 3b, and 6b) with 108 participants.

<sup>1407</sup> In most of these experiments, the interaction between self-reference and valence was  
<sup>1408</sup> significant (see results of each experiment in supplementary materials). In the  
<sup>1409</sup> mini-meta-analysis, we analyzed the valence effect for self-referential condition and  
<sup>1410</sup> other-referential condition separately.

<sup>1411</sup> For the self-referential condition, we found the same pattern as in the first part of  
<sup>1412</sup> results. That is we found significant differences between positive and neutral as well as  
<sup>1413</sup> positive and negative, but not neutral and negative. The effect size of RT between positive  
<sup>1414</sup> and negative is Cohen's  $d = -0.89 \pm 0.12$ , 95% CI [-1.11 -0.66]; on  $d'$  was Cohen's  $d = 0.61$   
<sup>1415</sup>  $\pm 0.09$ , 95% CI [0.44 0.78]. The effect was also observed between positive and neutral  
<sup>1416</sup> condition, RT: Cohen's  $d = -0.76 \pm 0.13$ , 95% CI [-1.01 -0.50];  $d'$ : Cohen's  $d = 0.69 \pm$   
<sup>1417</sup> 0.14, 95% CI [0.42 0.96]. And the difference between neutral and bad conditions are not  
<sup>1418</sup> significant, RT: Cohen's  $d = 0.03 \pm 0.13$ , 95% CI [-0.22 0.29];  $d'$ : Cohen's  $d = 0.08 \pm 0.08$ ,  
<sup>1419</sup> 95% CI [-0.07 0.24]. See Figure 31 the middle panel.

<sup>1420</sup> For the other-referential condition, we found that only the difference between positive  
<sup>1421</sup> and negative on RT was significant, all the other conditions were not. The effect size of RT  
<sup>1422</sup> between positive and negative is Cohen's  $d = -0.28 \pm 0.05$ , 95% CI [-0.38 -0.17]; on  $d'$  was  
<sup>1423</sup> Cohen's  $d = -0.02 \pm 0.08$ , 95% CI [-0.17 0.13]. The effect was not observed between

<sup>1424</sup> positive and neutral condition, RT: Cohen's  $d = -0.12 \pm 0.10$ , 95% CI [-0.31 0.06];  $d'$ :  
<sup>1425</sup> Cohen's  $d = 0.01 \pm 0.08$ , 95% CI [-0.16 0.17]. And the difference between neutral and bad  
<sup>1426</sup> conditions are not significant, RT: Cohen's  $d = 0.14 \pm 0.09$ , 95% CI [-0.03 0.31];  $d'$ :  
<sup>1427</sup> Cohen's  $d = 0.05 \pm 0.07$ , 95% CI [-0.08 0.18]. See Figure 31 right panel.

<sup>1428</sup> **Generalizability of the valence effect**

<sup>1429</sup> In this part, we reported the results from experiment 4 in which either moral valence  
<sup>1430</sup> or self-reference were manipulated as task-irrelevant stimuli.

<sup>1431</sup> For experiment 4a, when self-reference was the target and moral valence was  
<sup>1432</sup> task-irrelevant, we found that only under the implicit self-referential condition, i.e., when  
<sup>1433</sup> the moral words were presented as task irrelevant stimuli, there was the main effect of  
<sup>1434</sup> valence and interaction between valence and reference for both  $d$  prime and RT (See  
<sup>1435</sup> supplementary results for the detailed statistics). For  $d$  prime, we found good-self  
<sup>1436</sup> condition ( $2.55 \pm 0.86$ ) had higher  $d$  prime than bad-self condition ( $2.38 \pm 0.80$ ); good self  
<sup>1437</sup> condition was also higher than neutral self ( $2.45 \pm 0.78$ ) but there was not statistically  
<sup>1438</sup> significant, while the neutral-self condition was higher than bad self condition and not  
<sup>1439</sup> significant neither. For reaction times, good-self condition ( $654.26 \pm 67.09$ ) were faster  
<sup>1440</sup> relative to bad-self condition ( $665.64 \pm 64.59$ ), and over neutral-self condition ( $664.26 \pm$   
<sup>1441</sup>  $64.71$ ). The difference between neutral-self and bad-self conditions were not significant.  
<sup>1442</sup> However, for the other-referential condition, there was no significant differences between  
<sup>1443</sup> different valence conditions. See Figure 32.

<sup>1444</sup> For experiment 4b, when valence was the target and the identity was task-irrelevant,  
<sup>1445</sup> we found a strong valence effect (see supplementary results and Figure 33, Figure 34).

<sup>1446</sup> In this experiment, the advantage of good-self condition can only be disentangled by  
<sup>1447</sup> comparing the self-referential and other-referential conditions. Therefore, we calculated the  
<sup>1448</sup> differences between the valence effect under self-referential and other referential conditions

1449 and used the weighted variance as the variance of this differences. We found this  
1450 modulation effect on RT. The valence effect of RT was stronger in self-referential than  
1451 other-referential for the Good vs. Neutral condition ( $-0.33 \pm 0.01$ ), and to a less extent the  
1452 Good vs. Bad condition ( $-0.17 \pm 0.01$ ). While the size of the other effect's CI included  
1453 zero, suggesting those effects didn't differ from zero. See Figure 35.

1454 **Specificity of valence effect**

1455 In this part, we analyzed the results from experiment 5, which included positive,  
1456 neutral, and negative valence from four different domains: morality, emotion, aesthetics of  
1457 human, and aesthetics of scene. We found interaction between valence and domain for both  
1458  $d$  prime and RT (match trials). A common pattern appeared in all four domains: each  
1459 domain showed a binary results instead of gradient on both  $d$  prime and RT. For morality,  
1460 aesthetics of human, and aesthetics of scene, there was a positivity effect where the positive  
1461 conditions had advantages over both neutral (greater  $d$  prime and faster RT), while neutral  
1462 and negative conditions didn't differ from each other. But for the emotional stimuli, there  
1463 was a reversed negativity effect: positive and neutral conditions were not significantly  
1464 different from each other but both had advantage over negative conditions. See  
1465 supplementary materials for detailed statistics. Also note that the effect size in moral  
1466 domain is smaller than the aesthetic domains (beauty of people and beauty of scene). See  
1467 Figure 36.

1468 **Self-reported personal distance**

1469 See Figure 37.

1470 **Correlation analyses**

1471 The reliability of questionnaires can be found in (Liu et al., 2020). We calculated the  
1472 correlation between the data from behavioral task and the questionnaire data. First, we  
1473 calculated the score for each scale based on their structure and factor loading, instead of  
1474 sum score (McNeish & Wolf, 2020). Then, we used SEM to estimate the correlation  
1475 because it can include measurement model and statistical model in a unified framework.

1476 To make sure that what we found were not false positive, we used two method to  
1477 ensure the robustness of our analysis. first, we split the data into two half: the data with  
1478 self and without, then, we used the conditional random forest to find the robust correlation  
1479 in the exploratory data (with self reference) that can be replicated in the confirmatory data  
1480 (without the self reference). The robust correlation were then analyzed using SEM

1481 Instead of use the exploratory correlation analysis, we used a more principled way to  
1482 explore the correlation between parameter of HDDM ( $v$ ,  $t$ , and  $a$ ) and scale scores and  
1483 person distance.

1484 We didn't find the correlation between scale scores and the parameters of HDDM,  
1485 but found weak correlation between personal distance and the parameter estimated from  
1486 Good and neutral conditions.

1487 First, boundary separation ( $a$ ) of moral good condition was correlated with both  
1488 Self-Bad distance ( $r = 0.198$ , 95% CI  $[], p = 0.0063$ ) and Neutral-Bad distance  
1489 ( $r = 0.1571$ , 95% CI  $[], p = 0.031$ ). At the same time, the non-decision time is negatively  
1490 correlated with Self-Bad distance ( $r = 0.169$ , 95% CI  $[], p = 0.0197$ ). See Figure 38.

1491 Second, we found the boundary separation of neutral condition is positively  
1492 correlated with the personal distance between self and good distance ( $r = 0.189$ , 95% CI  $[],$   
1493  $p = 0.036$ ), but negatively correlated with self-neutral distance( $r = -0.183$ , 95% CI  $[],$   
1494  $p = 0.042$ ). Also, the drift rate of the neutral condition is positively correlated with the  
1495 Self-Bad distance ( $r = 0.177$ , 95% CI  $[], p = 0.048$ ).a. See figure 39

1496 We also explored the correlation between behavioral data and questionnaire scores  
1497 separately for experiments with and without self-referential, however, the sample size is  
1498 very low for some conditions.

1499 **Discussion**

1500 **References**

- 1501 Anderson, E., Siegel, E. H., Bliss-Moreau, E., & Barrett, L. F. (2011). The visual impact  
1502 of gossip. *Science*, 332(6036), 1446–1448. <https://doi.org/10.1126/science.1201574>
- 1503 Brady, W. J., Gantman, A. P., & Van Bavel, J. J. (2020). Attentional capture helps  
1504 explain why moral and emotional content go viral. *Journal of Experimental  
1505 Psychology: General*, 149(4), 746–756. <https://doi.org/10.1037/xge0000673>
- 1506 Brainard, D. H. (1997). The psychophysics toolbox. *Spatial Vision*, 10(4), 433–436.  
1507 Journal Article.
- 1508 Bruner, J. S. (1957). On perceptual readiness. *Psychological Review*, 64(2), 123–152.  
1509 <https://doi.org/10.1037/h0043805>
- 1510 Bürkner, P.-C. (2017). Brms: An r package for bayesian multilevel models using stan.  
1511 *Journal of Statistical Software; Vol 1, Issue 1 (2017)*. Journal Article. Retrieved  
1512 from  
1513 <https://www.jstatsoft.org/v080/i01%0Ahttp://dx.doi.org/10.18637/jss.v080.i01>
- 1514 Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M., ...  
1515 Riddell, A. (2017). Stan: A probabilistic programming language. *Journal of  
1516 Statistical Software*, 76(1). Journal Article. <https://doi.org/10.18637/jss.v076.i01>
- 1517 Cooper, H., Hedges, L. V., & Valentine, J. C. (2009). *The handbook of research synthesis  
1518 and meta-analysis* (2nd ed.). Book, New York: Sage.

- 1519 DeCarlo, L. T. (1998). Signal detection theory and generalized linear models. *Psychological  
1520 Methods*, 3(2), 186–205. Journal Article. <https://doi.org/10.1037/1082-989X.3.2.186>
- 1521 DeScioli, P. (2016). The side-taking hypothesis for moral judgment. *Current Opinion in  
1522 Psychology*, 7, 23–27. <https://doi.org/10.1016/j.copsyc.2015.07.002>
- 1523 Faul, F., Erdfelder, E., Buchner, A., & Lang, A.-G. (2009). Statistical power analyses using  
1524 g\*Power 3.1: Tests for correlation and regression analyses. *Behavior Research  
1525 Methods*, 41(4), 1149–1160. <https://doi.org/10.3758/BRM.41.4.1149>
- 1526 Freitas, J. D., Cikara, M., Grossmann, I., & Schlegel, R. (2017). Origins of the belief in  
1527 good true selves. *Trends in Cognitive Sciences*, 21(9), 634–636.  
1528 <https://doi.org/10.1016/j.tics.2017.05.009>
- 1529 Gantman, A. P., & Van Bavel, J. J. (2014). The moral pop-out effect: Enhanced  
1530 perceptual awareness of morally relevant stimuli. *Cognition*, 132(1), 22–29.  
1531 <https://doi.org/10.1016/j.cognition.2014.02.007>
- 1532 Gantman, A. P., & Van Bavel, J. J. (2016). Exposure to justice diminishes moral  
1533 perception. *Journal of Experimental Psychology: General*, 145(12), 1728–1739.  
1534 <https://doi.org/10.1037/xge0000241>
- 1535 Gilbert, C. D., & Li, W. (2013). Top-down influences on visual processing. *Nature Reviews  
1536 Neuroscience*, 14(5), 350–363. <https://doi.org/10.1038/nrn3476>
- 1537 Goh, J. X., Hall, J. A., & Rosenthal, R. (2016). Mini meta-analysis of your own studies:  
1538 Some arguments on why and a primer on how. *Social and Personality Psychology  
1539 Compass*, 10(10), 535–549. Journal Article. <https://doi.org/10.1111/spc3.12267>
- 1540 Hu, C.-P., Lan, Y., Macrae, C. N., & Sui, J. (2020). Good me bad me: Does valence  
1541 influence self-prioritization during perceptual decision-making? *Collabra:  
1542 Psychology*, 6(1), 20. Journal Article. <https://doi.org/10.1525/collabra.301>
- 1543 Juechems, K., & Summerfield, C. (2019). Where does value come from? *Trends in*

- 1544        *Cognitive Sciences*, 23(10), 836–850. <https://doi.org/10.1016/j.tics.2019.07.012>
- 1545    Kirby, K. N., & Gerlanc, D. (2013). BootES: An r package for bootstrap confidence  
1546        intervals on effect sizes. *Behavior Research Methods*, 45(4), 905–927.  
1547        <https://doi.org/10.3758/s13428-013-0330-5>
- 1548    Liu, Q., Wang, F., Yan, W., Peng, K., Sui, J., & Hu, C.-P. (2020). Questionnaire data from  
1549        the revision of a chinese version of free will and determinism plus scale. *Journal of*  
1550        *Open Psychology Data*, 8(1), 1. Journal Article. <https://doi.org/10.5334/jopd.49/>
- 1551    Matzke, D., & Wagenmakers, E.-J. (2009). Psychological interpretation of the ex-gaussian  
1552        and shifted wald parameters: A diffusion model analysis. *Psychonomic Bulletin &*  
1553        *Review*, 16(5), 798–817. <https://doi.org/10.3758/PBR.16.5.798>
- 1554    McNeish, D., & Wolf, M. G. (2020). Thinking twice about sum scores. *Behavior Research*  
1555        *Methods*. <https://doi.org/10.3758/s13428-020-01398-0>
- 1556    Pelli, D. G. (1997). The videotoolbox software for visual psychophysics: Transforming  
1557        numbers into movies. *Spatial Vision*, 10(4), 437–442. Journal Article.
- 1558    Rouder, J. N., & Lu, J. (2005). An introduction to bayesian hierarchical models with an  
1559        application in the theory of signal detection. *Psychonomic Bulletin & Review*,  
1560        12(4), 573–604. Journal Article. <https://doi.org/10.3758/bf03196750>
- 1561    Rousselet, G. A., & Wilcox, R. R. (2019). Reaction times and other skewed distributions:  
1562        Problems with the mean and the median. *Meta-Psychology*. preprint.  
1563        <https://doi.org/10.1101/383935>
- 1564    Shadlen, M. N., & Kiani, R. (2013). Decision making as a window on cognition. *Neuron*,  
1565        80(3), 791–806. <https://doi.org/10.1016/j.neuron.2013.10.047>
- 1566    Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2013). Life after p-hacking. Conference  
1567        Proceedings. <https://doi.org/10.2139/ssrn.2205186>
- 1568    Stanislaw, H., & Todorov, N. (1999). Calculation of signal detection theory measures.

- 1569        *Behavior Research Methods, Instruments, & Computers*, 31(1), 137–149. Journal  
1570        Article. <https://doi.org/10.3758/BF03207704>
- 1571        Stolier, R. M., & Freeman, J. B. (2016). Functional and temporal considerations for  
1572        top-down influences in social perception. *Psychological Inquiry*, 27(4), 352–357.  
1573        <https://doi.org/10.1080/1047840X.2016.1216034>
- 1574        Strohminger, N., Knobe, J., & Newman, G. (2017). The true self: A psychological concept  
1575        distinct from the self: *Perspectives on Psychological Science*.  
1576        <https://doi.org/10.1177/1745691616689495>
- 1577        Sui, J., He, X., & Humphreys, G. W. (2012). Perceptual effects of social salience: Evidence  
1578        from self-prioritization effects on perceptual matching. *Journal of Experimental  
1579        Psychology: Human Perception and Performance*, 38(5), 1105–1117. Journal  
1580        Article. <https://doi.org/10.1037/a0029792>
- 1581        Wiecki, T. V., Sofer, I., & Frank, M. J. (2013). HDDM: Hierarchical bayesian estimation of  
1582        the drift-diffusion model in python. *Frontiers in Neuroinformatics*, 7.  
1583        <https://doi.org/10.3389/fninf.2013.00014>
- 1584        Xiao, Y. J., Coppin, G., & Bavel, J. J. V. (2016). Perceiving the world through  
1585        group-colored glasses: A perceptual model of intergroup relations. *Psychological  
1586        Inquiry*, 27(4), 255–274. <https://doi.org/10.1080/1047840X.2016.1199221>
- 1587        Zhang, H., Chen, K., Schlegel, R., Hicks, J., & Chen, C. (2019). The authentic moral self:  
1588        Dynamic interplay between perceived authenticity and moral behaviors in the  
1589        workplace. *Collabra: Psychology*, 5(1), 48. <https://doi.org/10.1525/collabra.260>

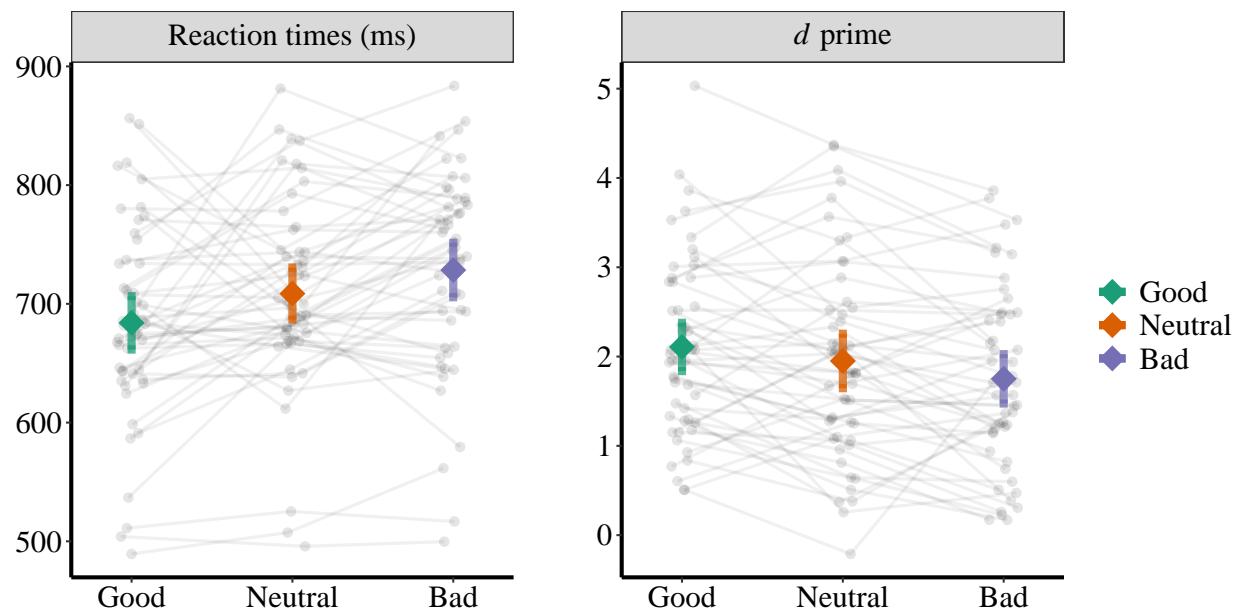


Figure 1. RT and  $d$  prime of Experiment 1a.

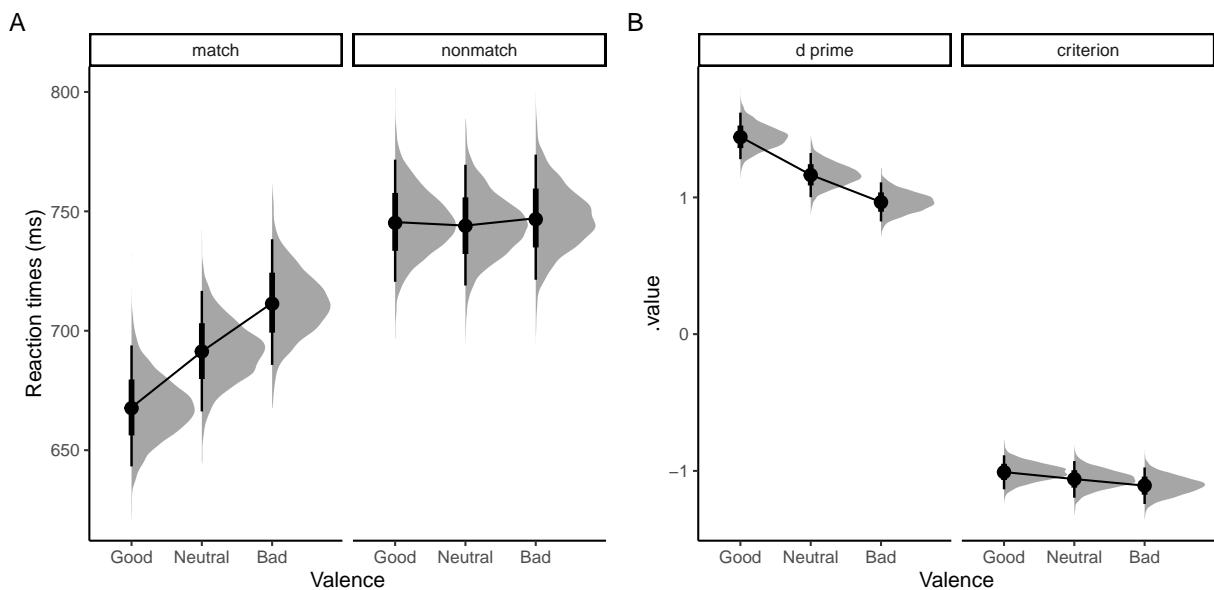


Figure 2. Exp1a: Results of Bayesian GLM analysis.

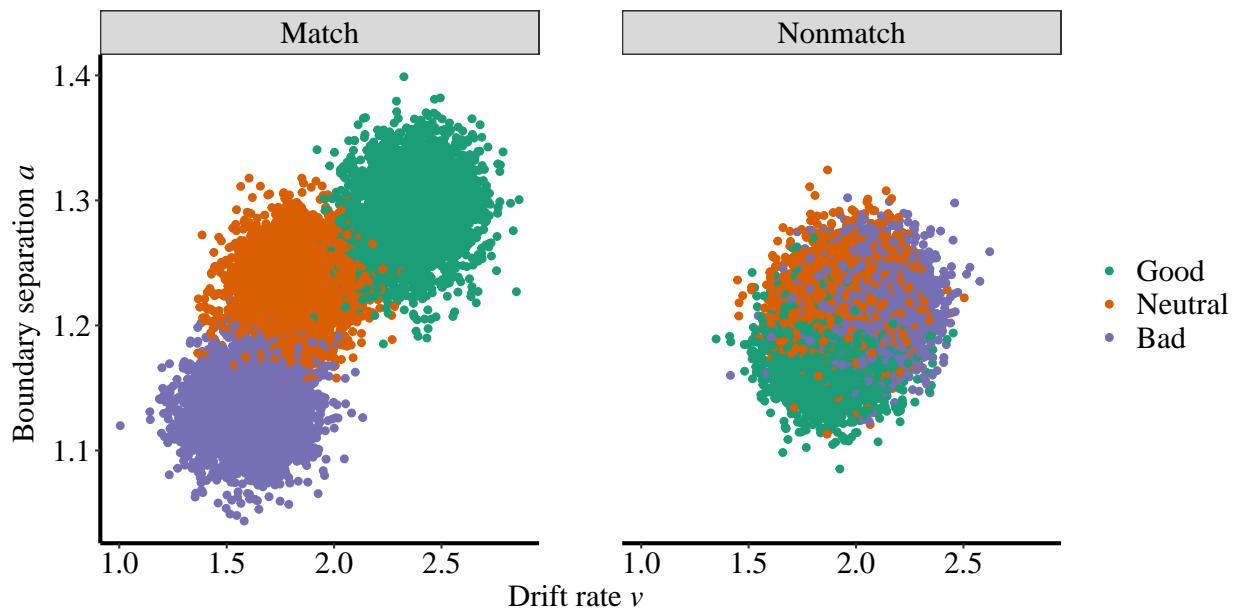


Figure 3. Exp1a: Results of HDDM.

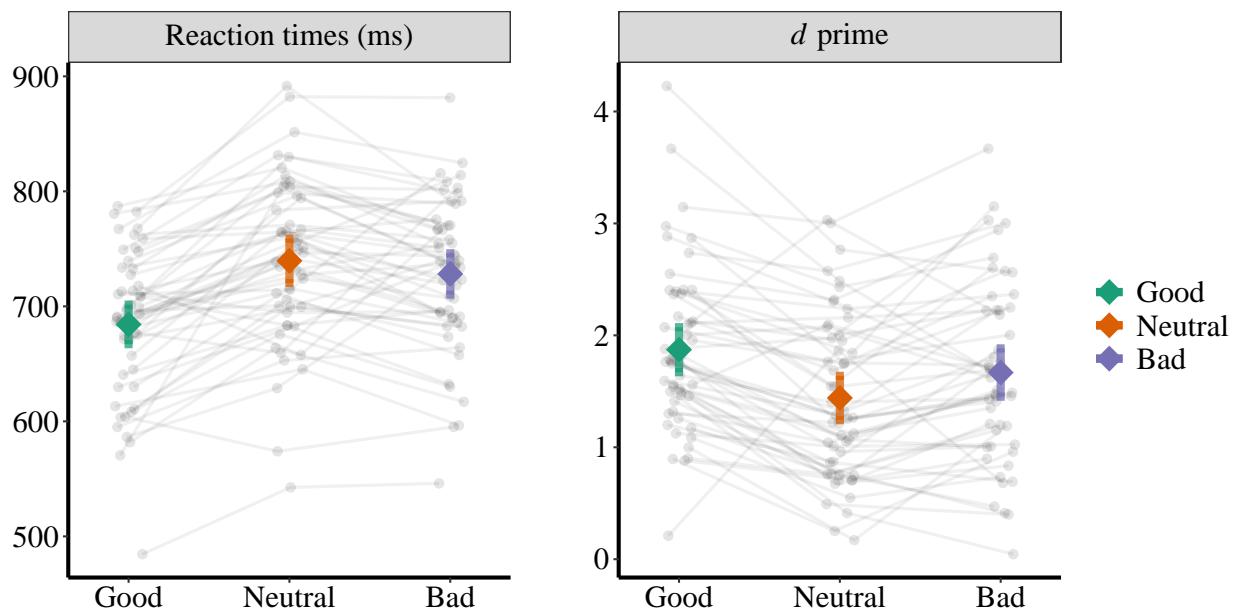


Figure 4. RT and  $d'$  of Experiment 1b.

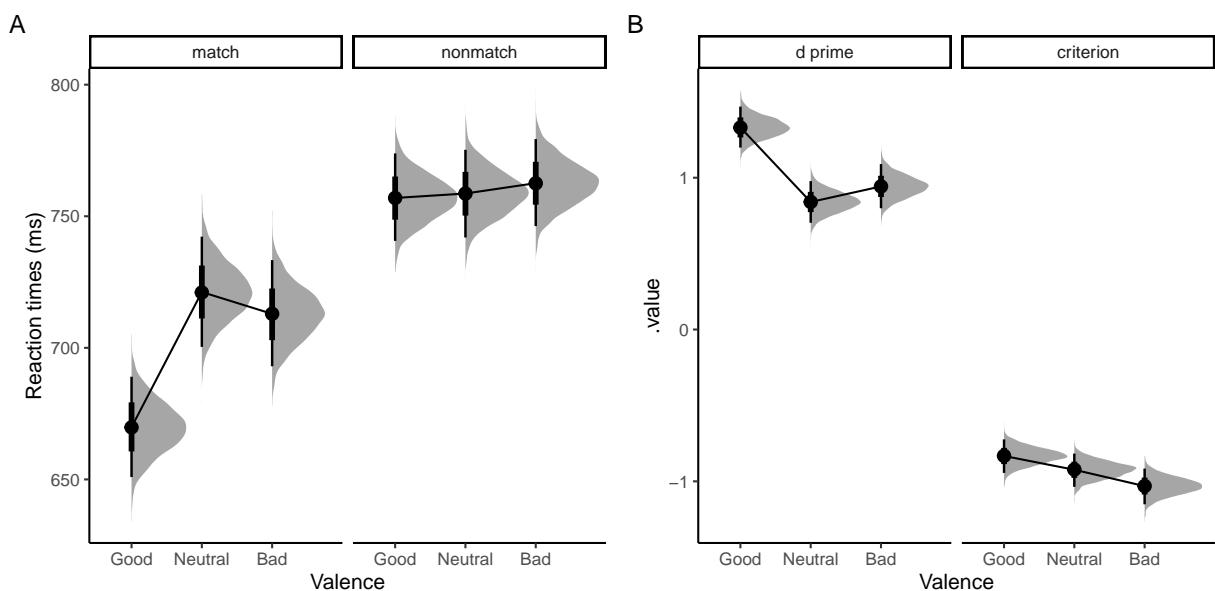


Figure 5. Exp1b: Results of Bayesian GLM analysis.

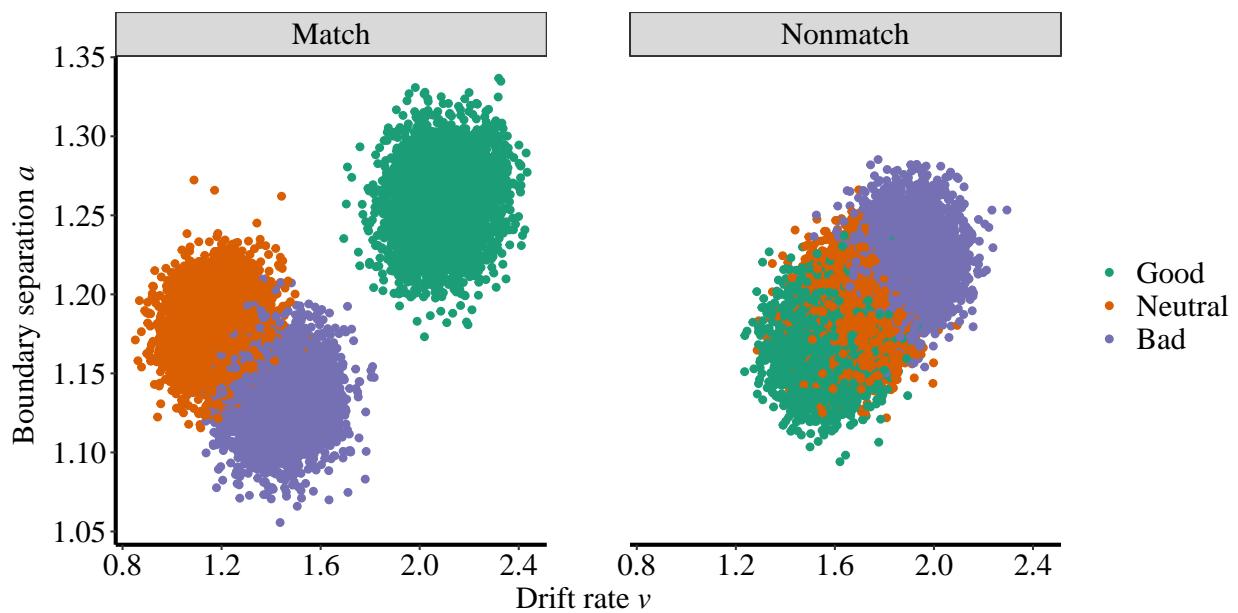


Figure 6. Exp1b: Results of HDDM.

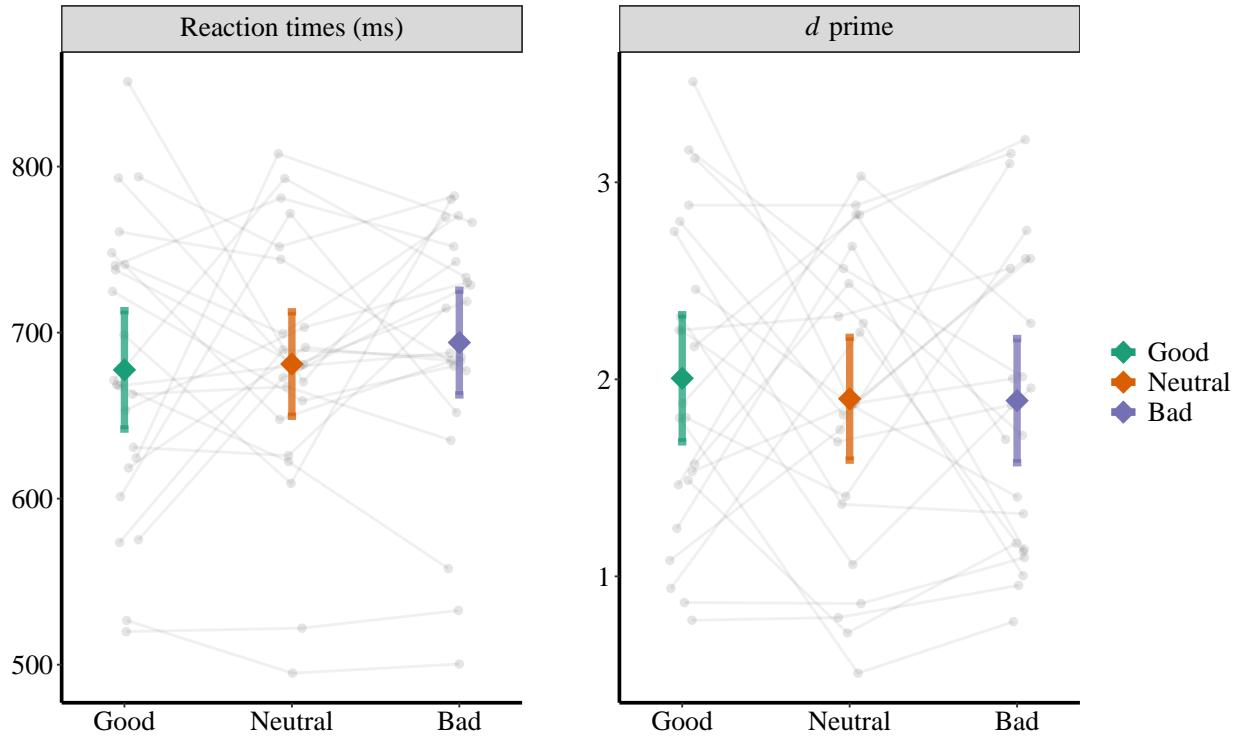


Figure 7. RT and  $d'$  prime of Experiment 1c.

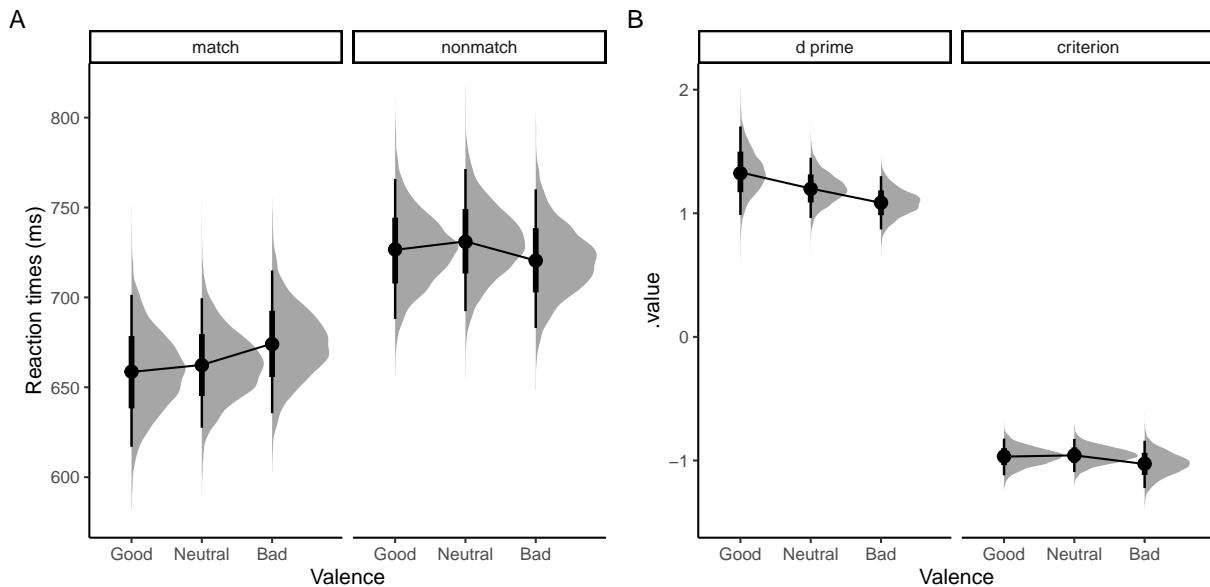


Figure 8. Exp1c: Results of Bayesian GLM analysis.

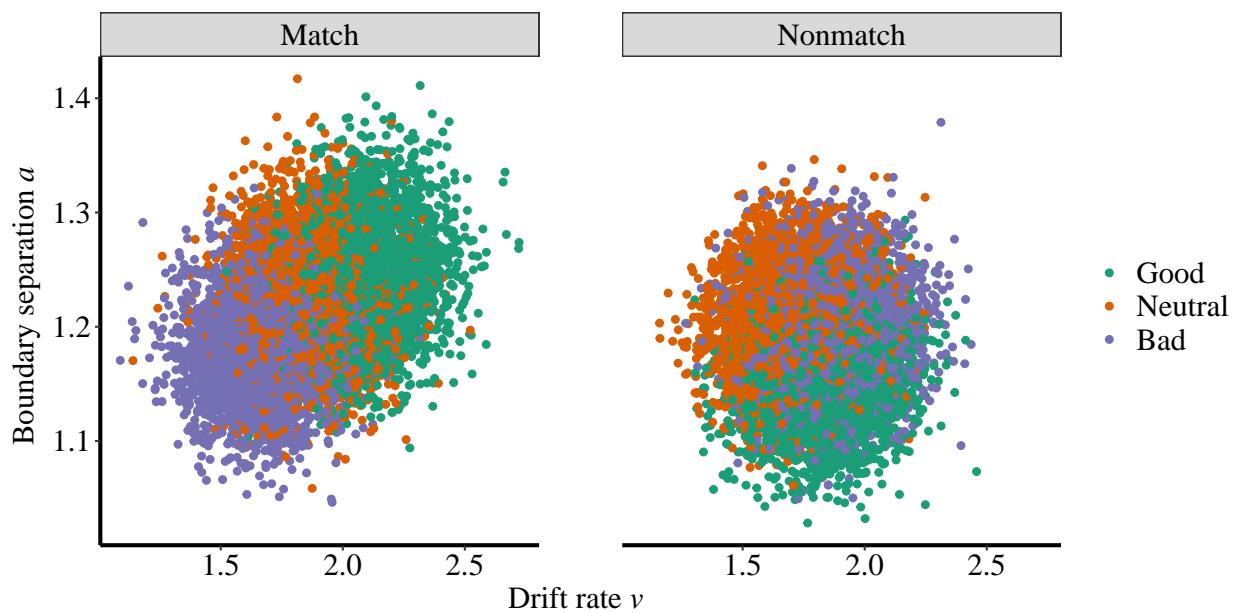


Figure 9. Exp1c: Results of HDDM.

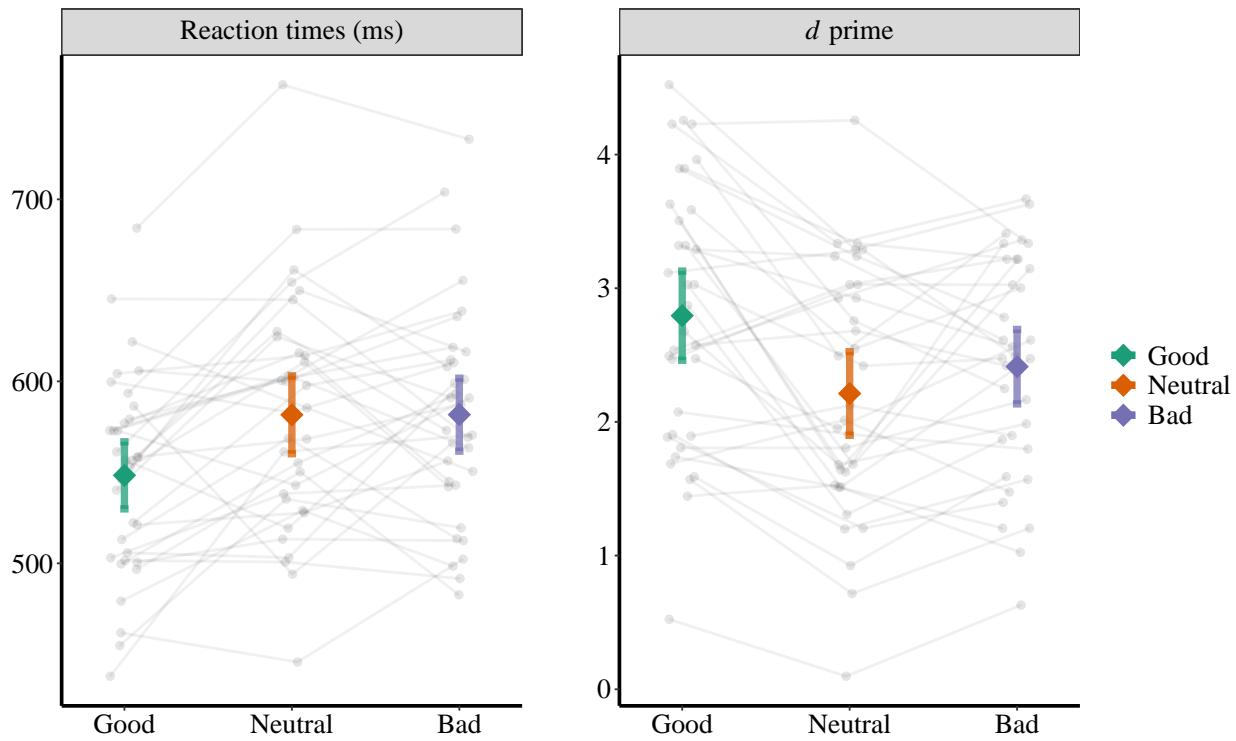


Figure 10. RT and  $d'$  of Experiment 2.

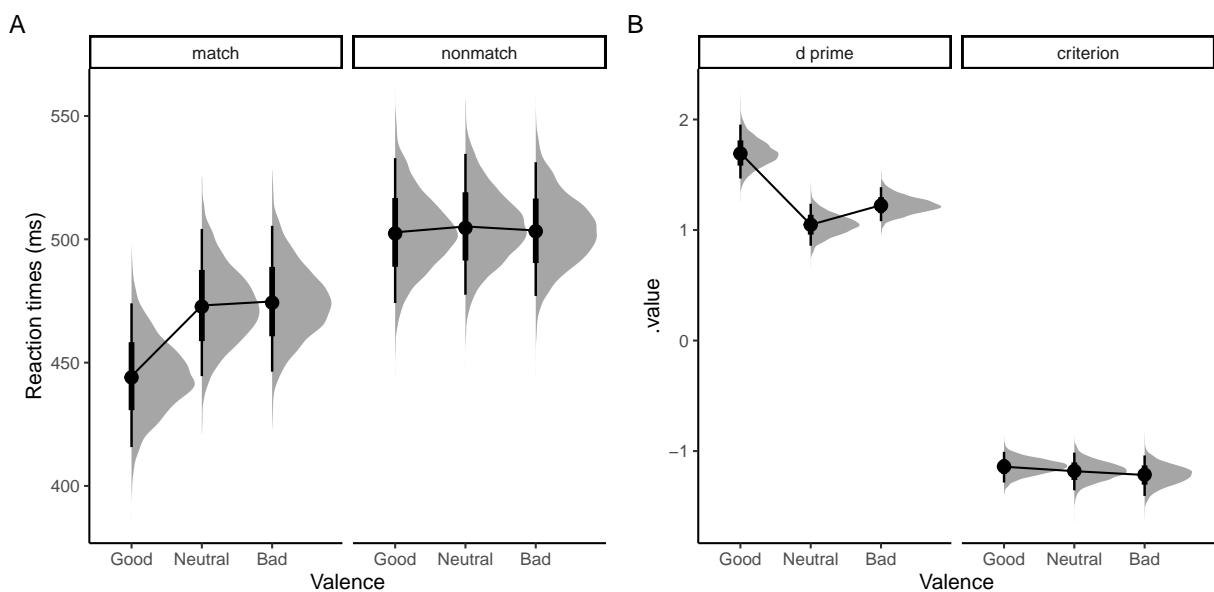


Figure 11. Exp2: Results of Bayesian GLM analysis.

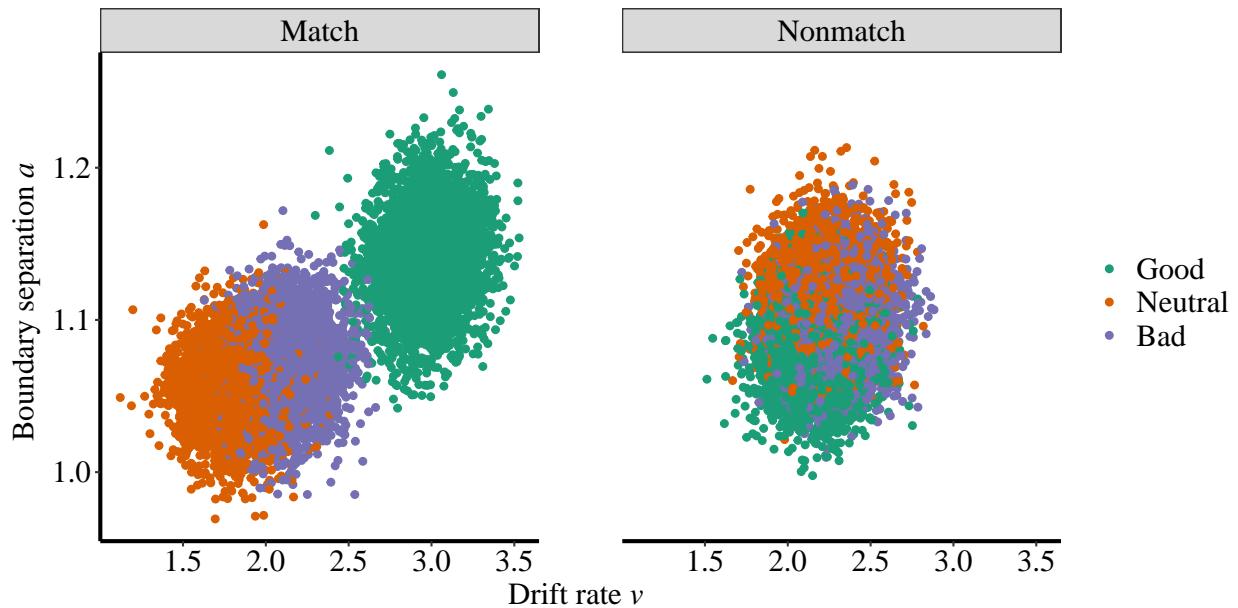


Figure 12. Exp2: Results of HDDM.

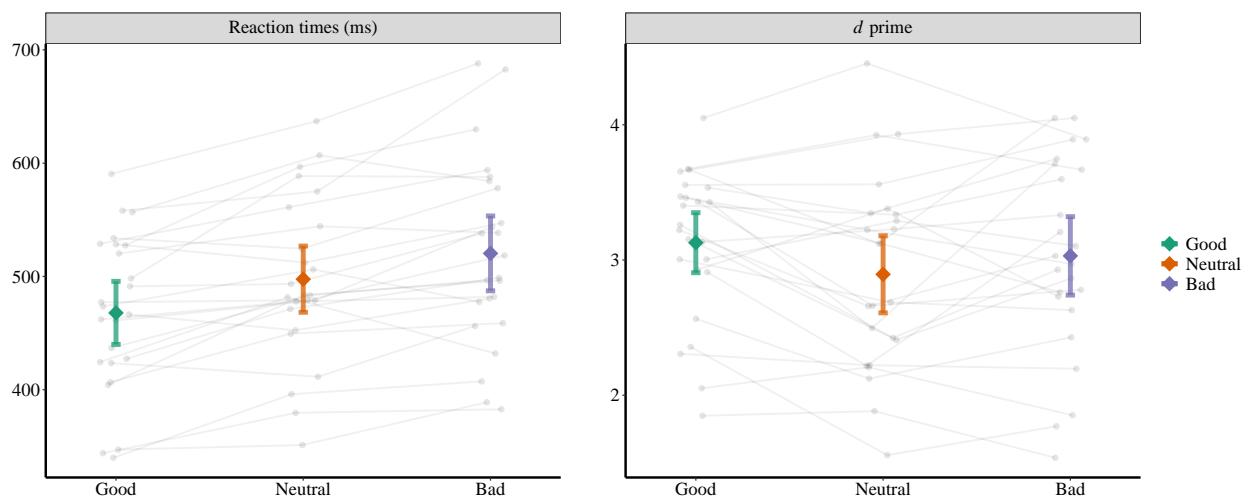


Figure 13. RT and  $d'$  prime of Experiment 6a.

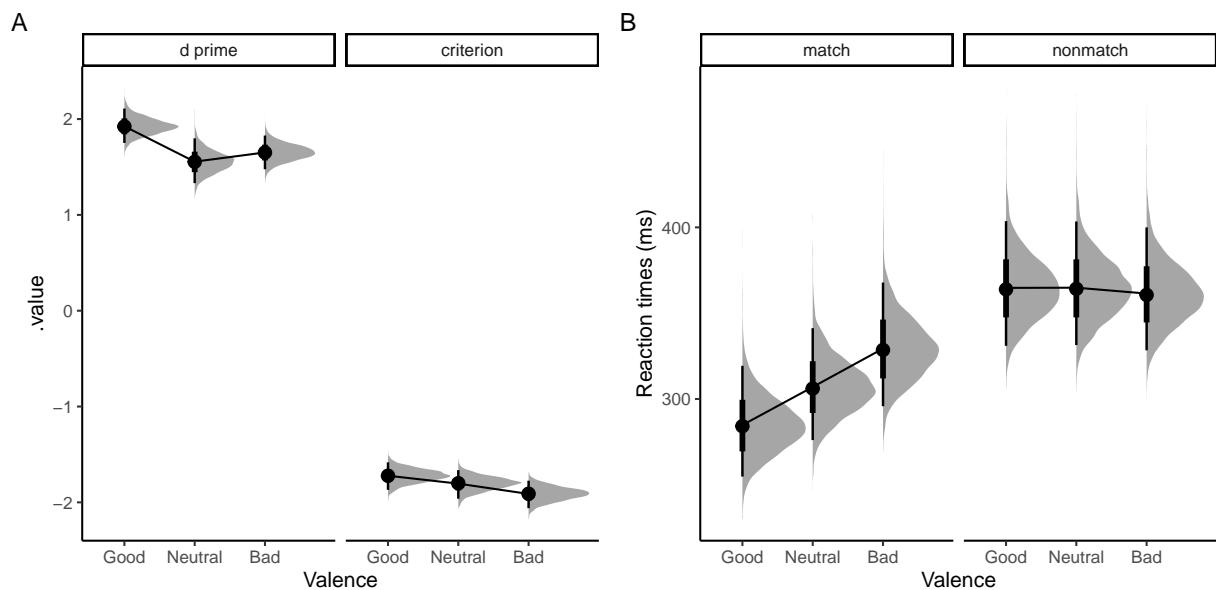


Figure 14. Exp6a: Results of Bayesian GLM analysis.

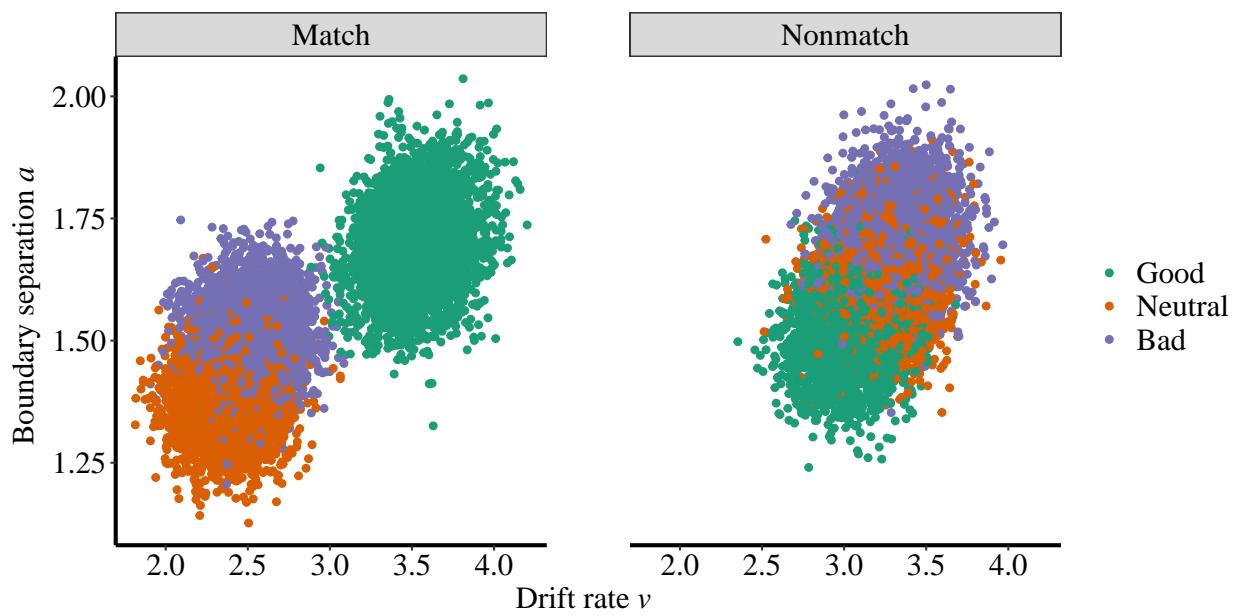


Figure 15. exp6a: Results of HDDM.

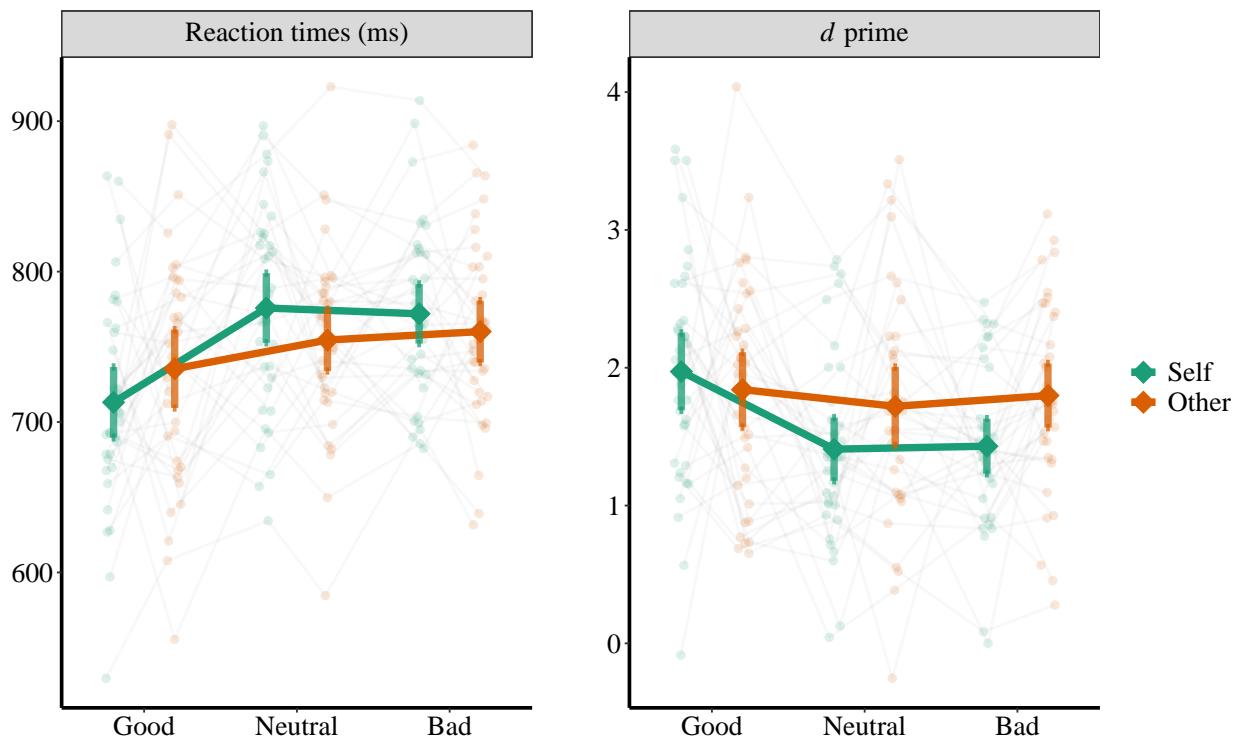


Figure 16. RT and  $d'$  prime of Experiment 3a.

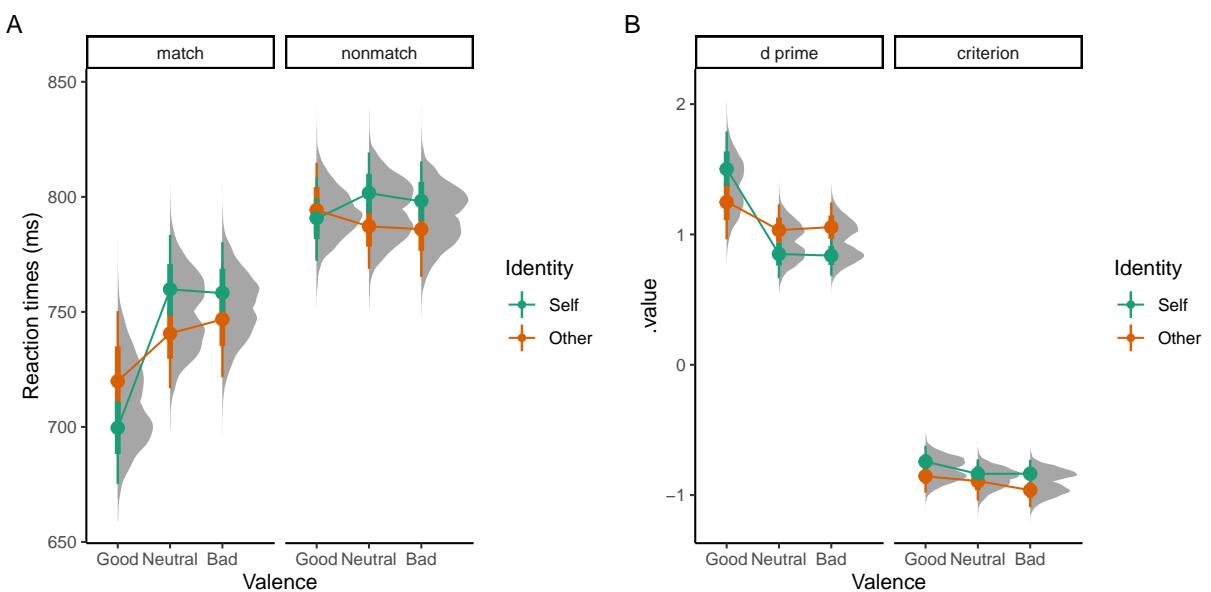


Figure 17. Exp3a: Results of Bayesian GLM analysis.

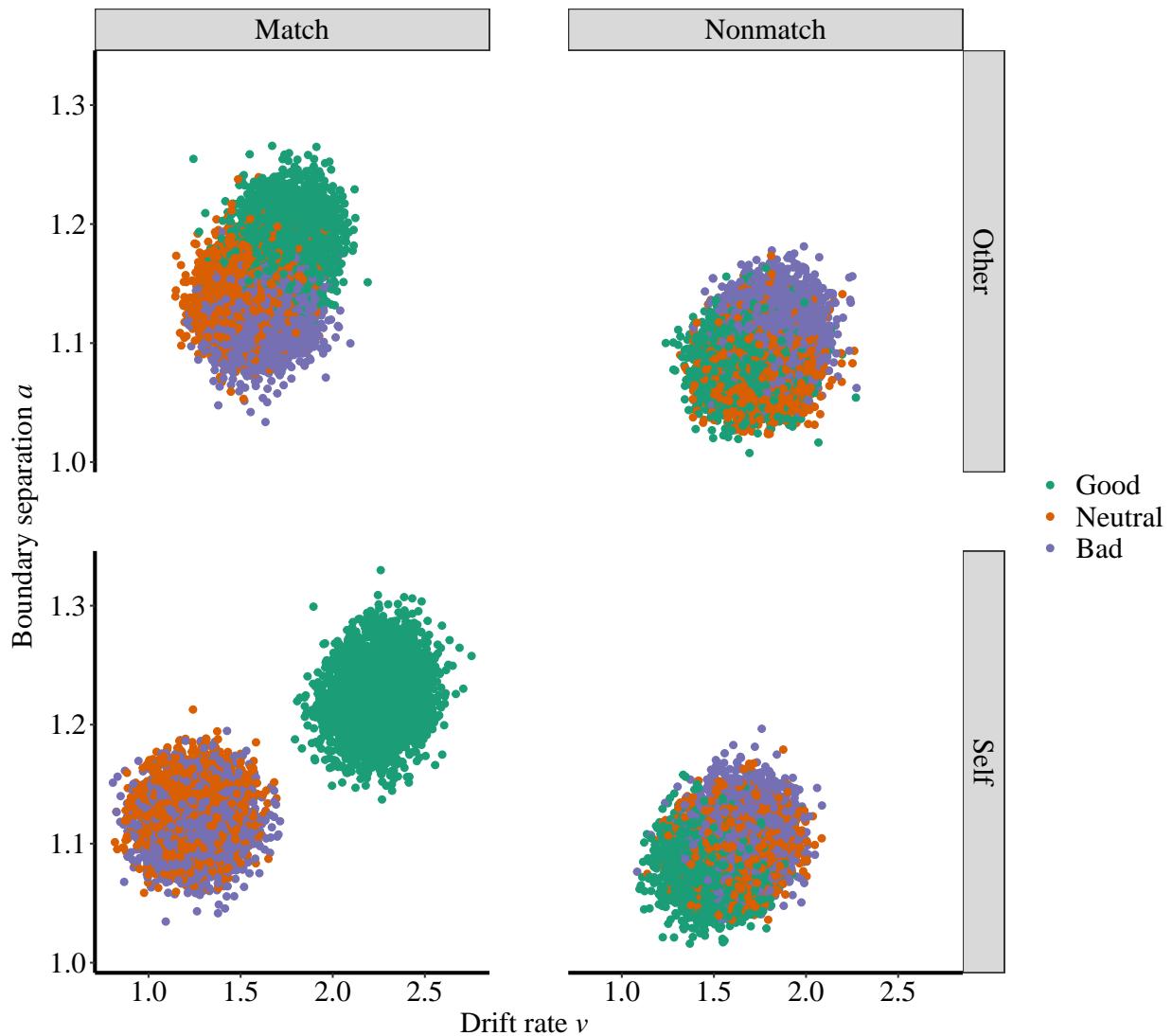


Figure 18. Exp3a: Results of HDDM.

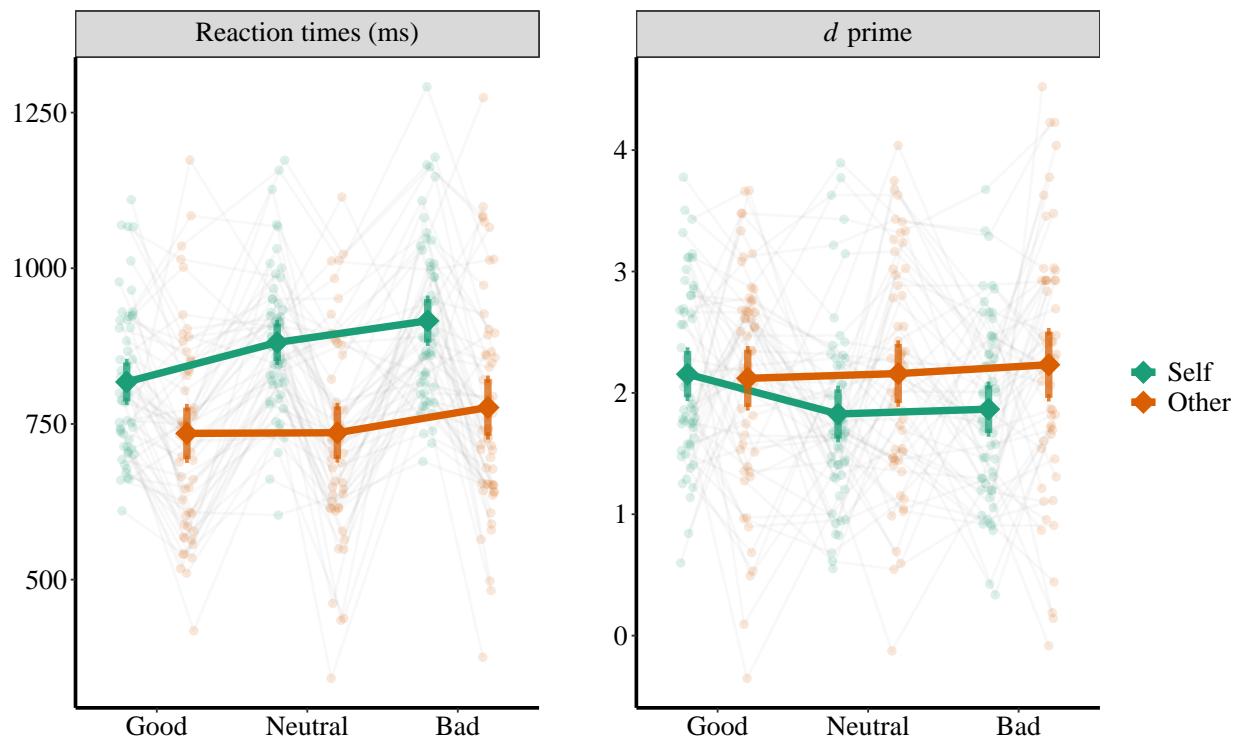


Figure 19. RT and  $d$  prime of Experiment 3b.

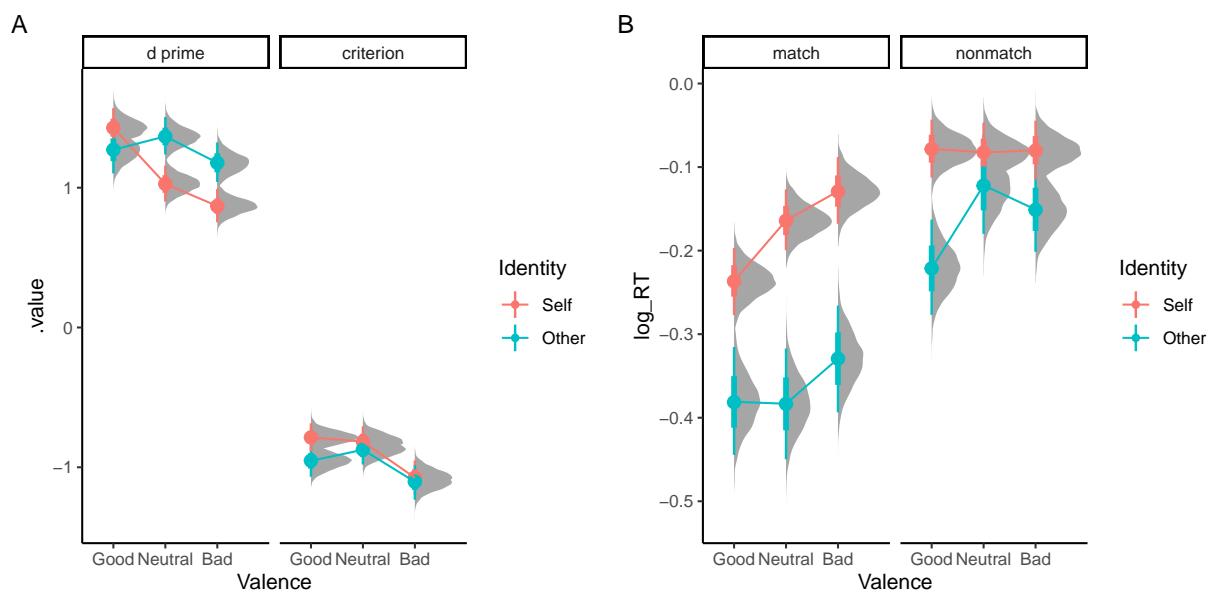


Figure 20. exp3b: Results of Bayesian GLM analysis.

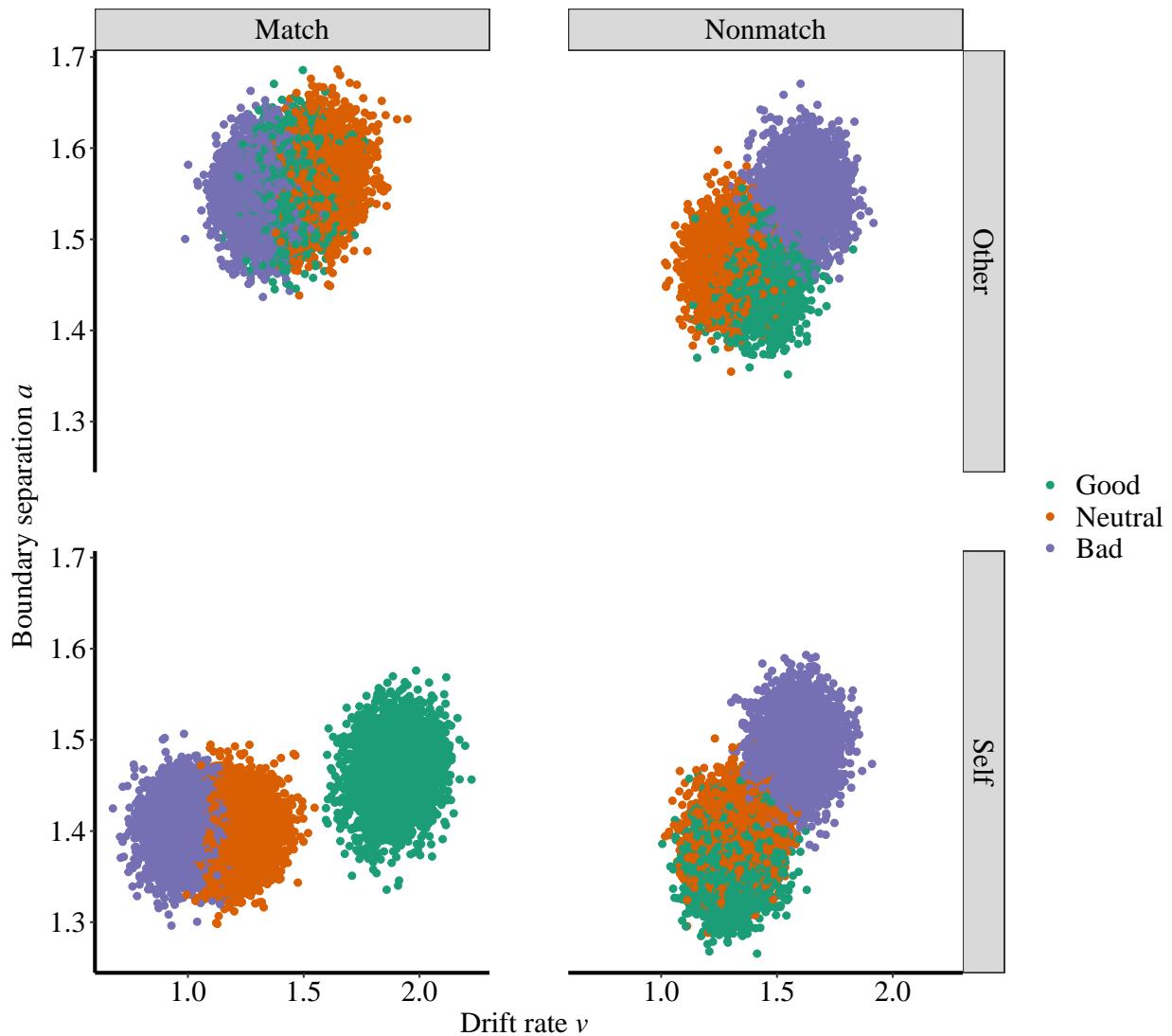


Figure 21. exp3b: Results of HDDM.

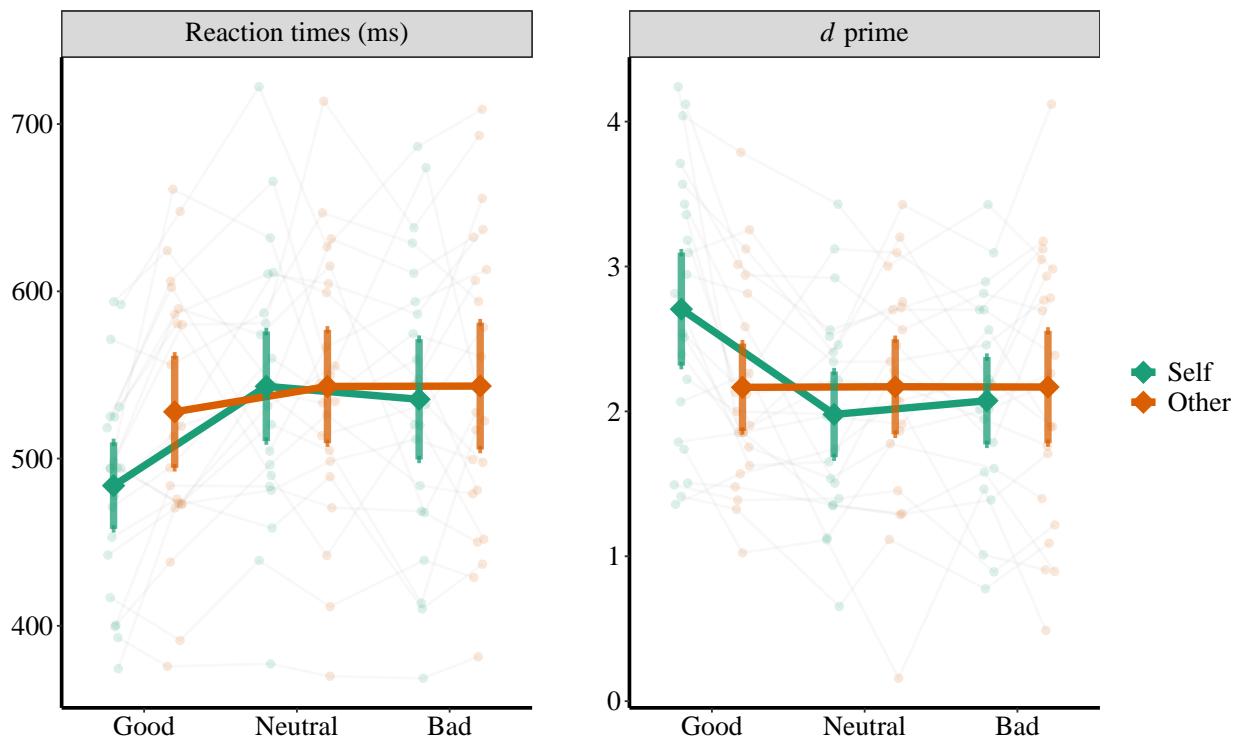


Figure 22. RT and  $d$  prime of Experiment 6b.

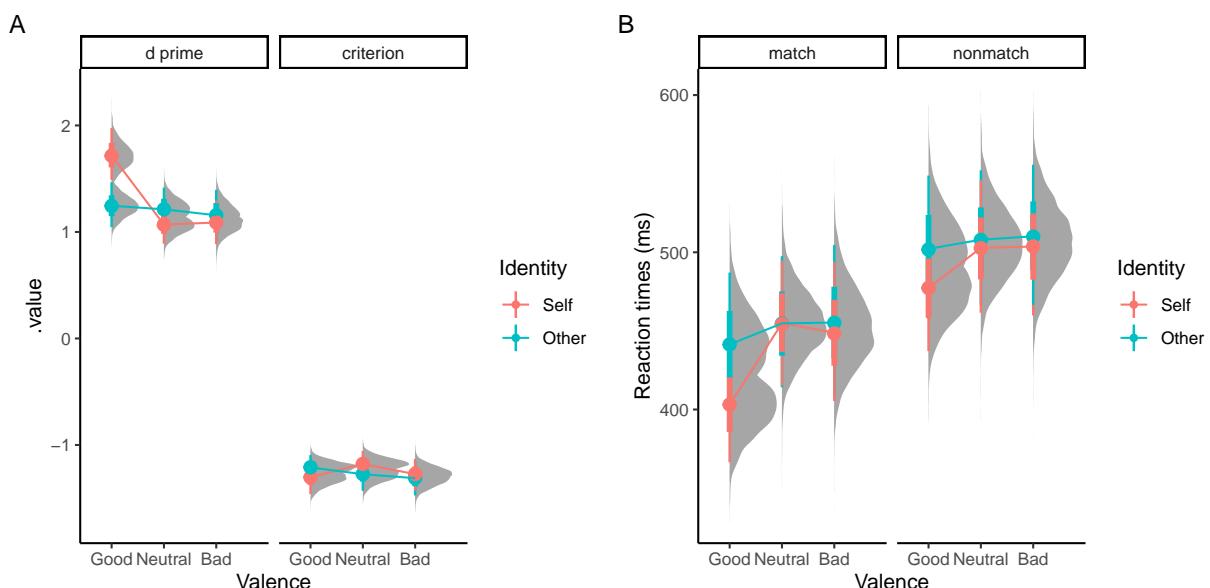


Figure 23. exp6b\_d1: Results of Bayesian GLM analysis.

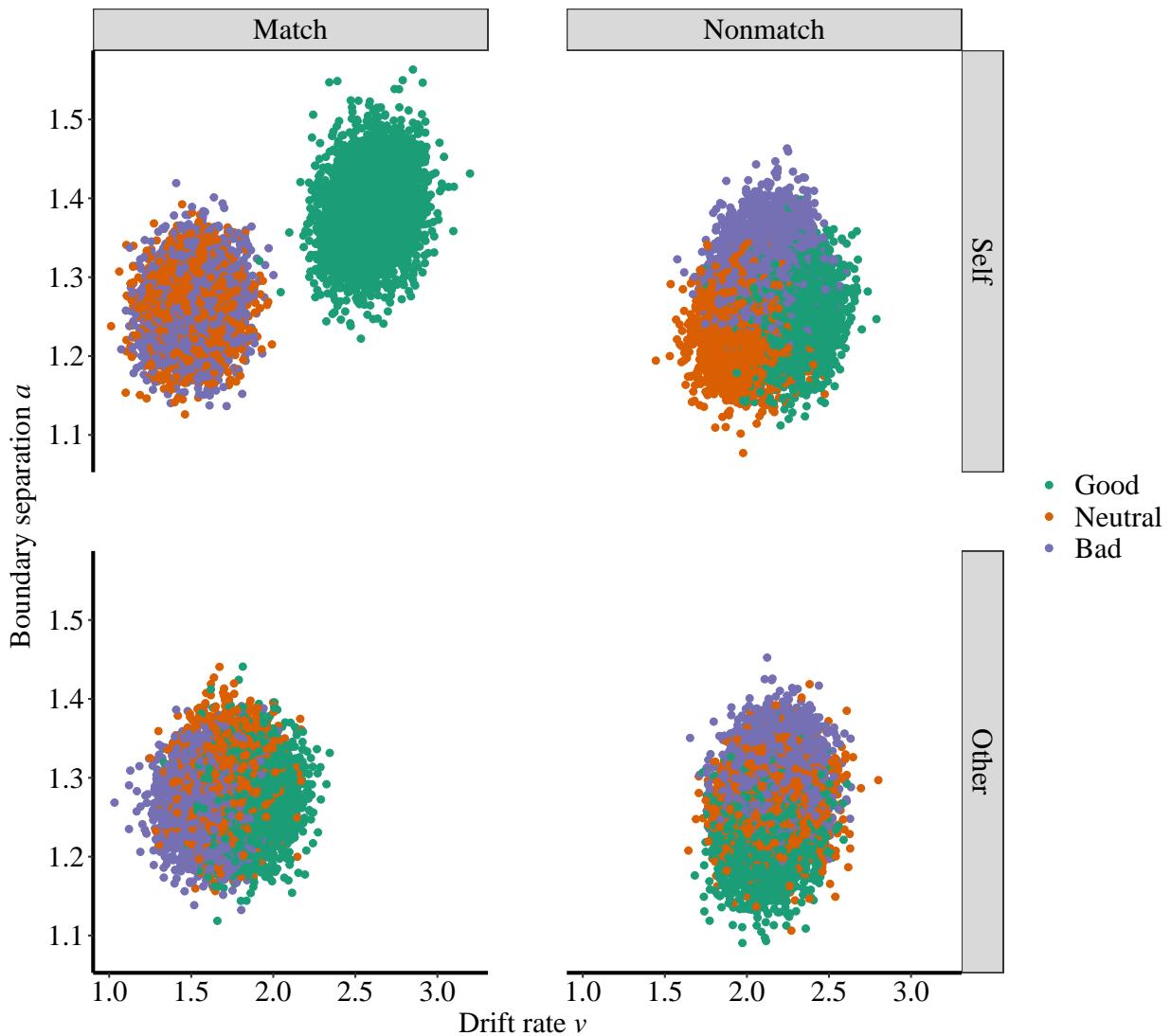


Figure 24. exp6b: Results of HDDM (Day 1).

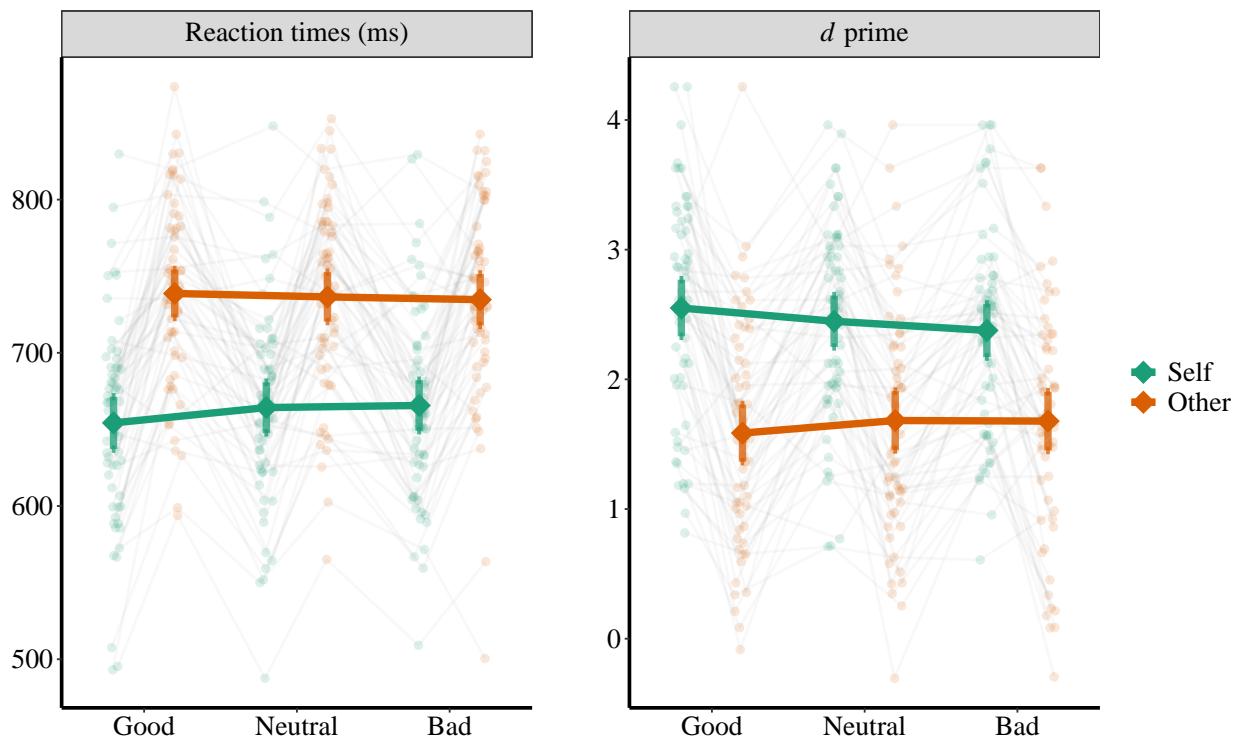


Figure 25. RT and  $d'$  of Experiment 4a.

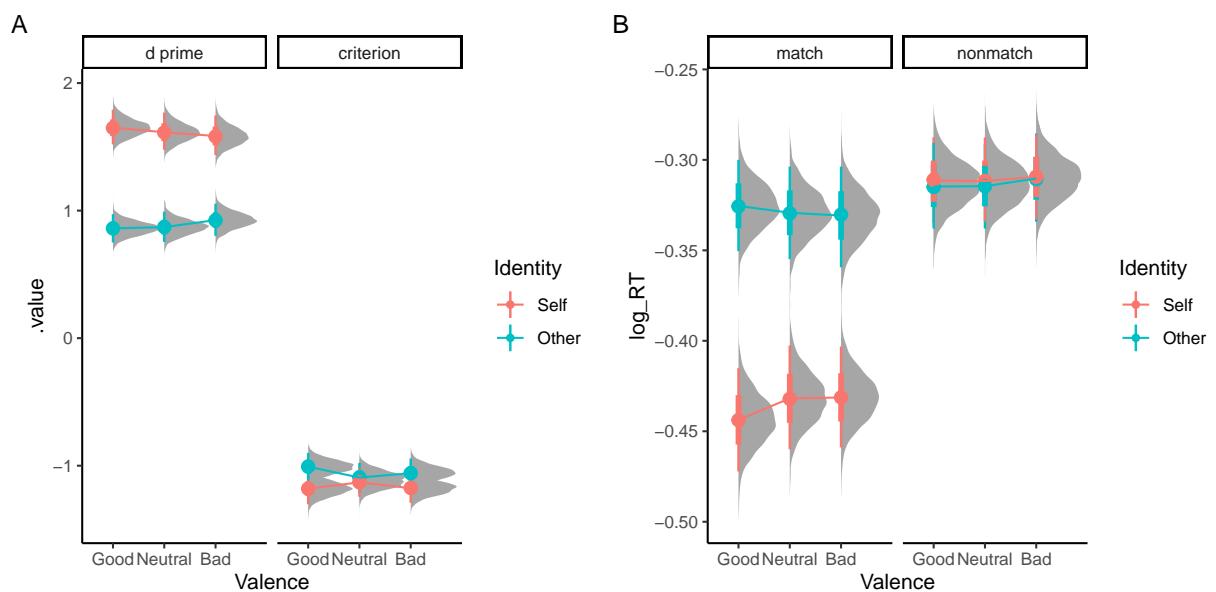


Figure 26. exp4a: Results of Bayesian GLM analysis.

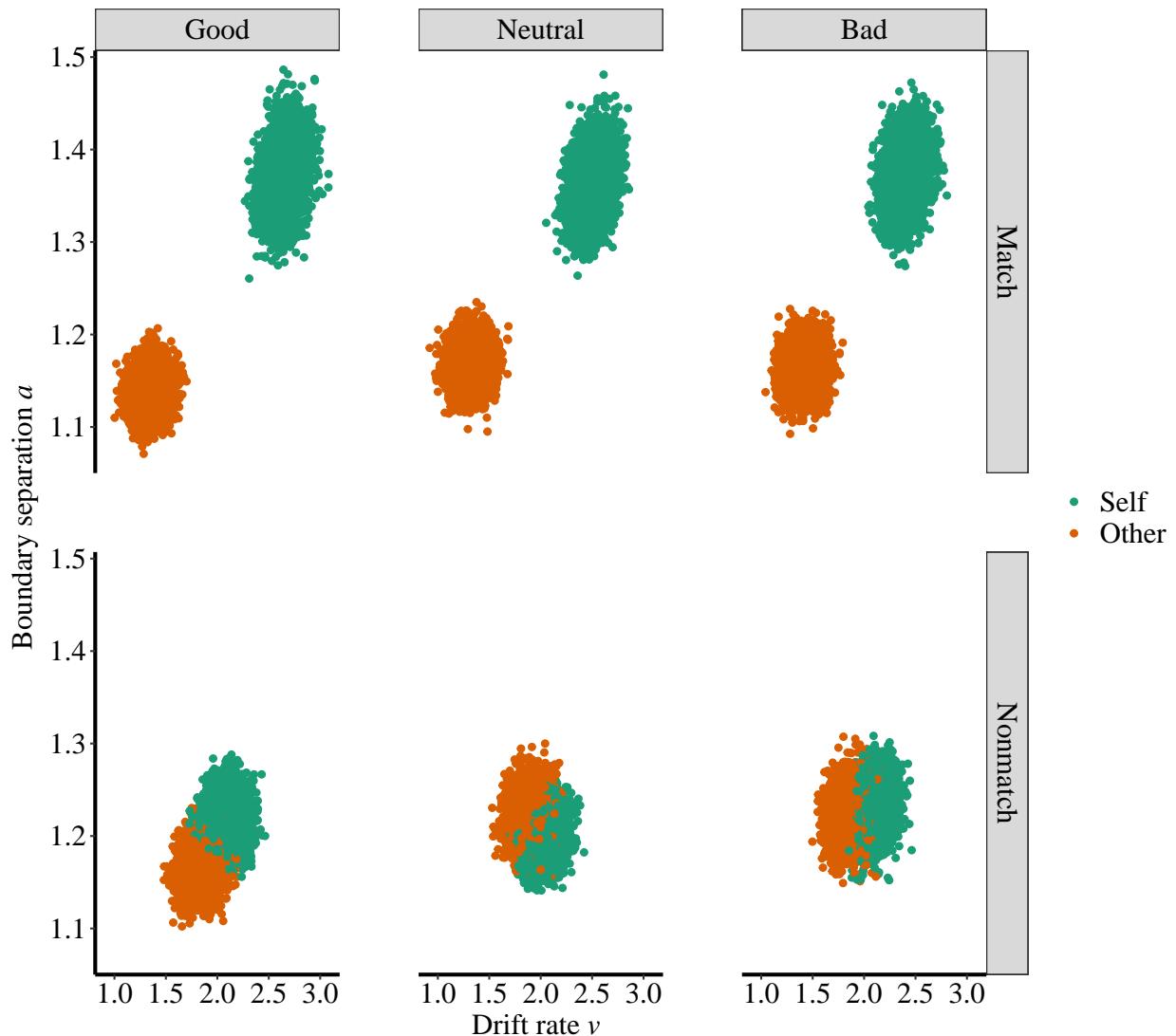


Figure 27. exp4a: Results of HDDM.

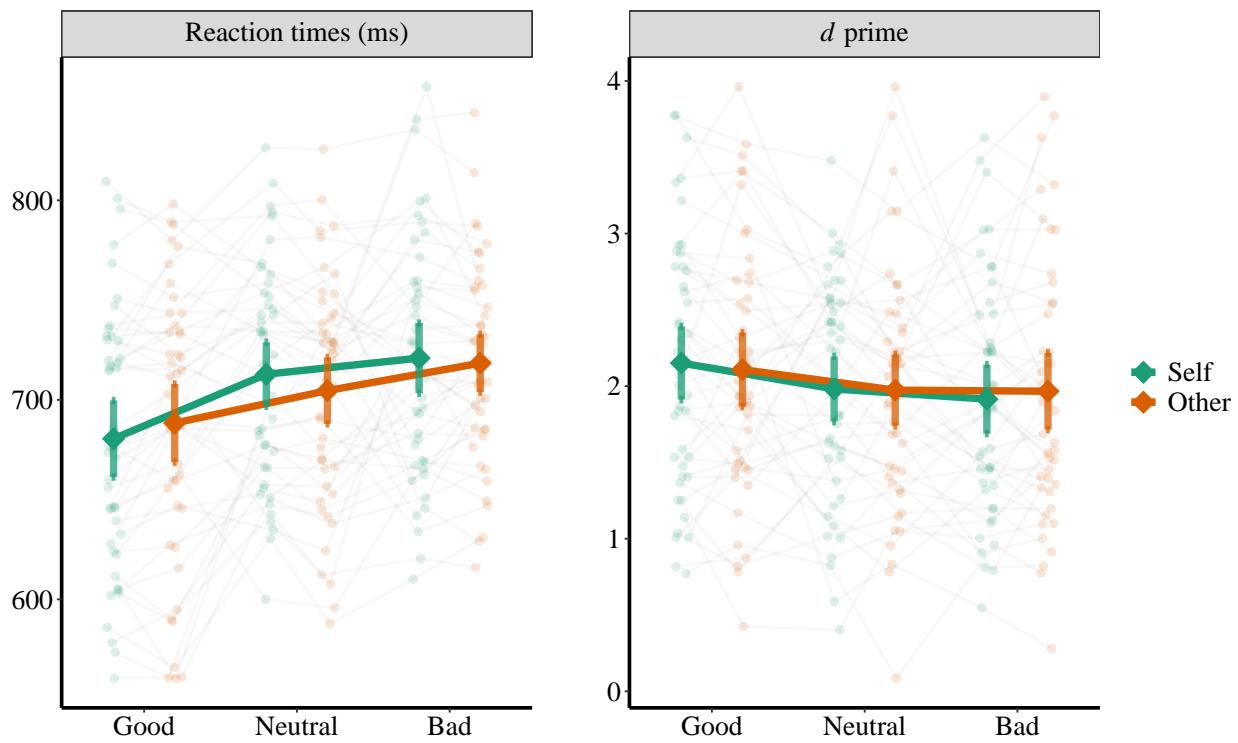


Figure 28. RT and  $d'$  prime of Experiment 4b.

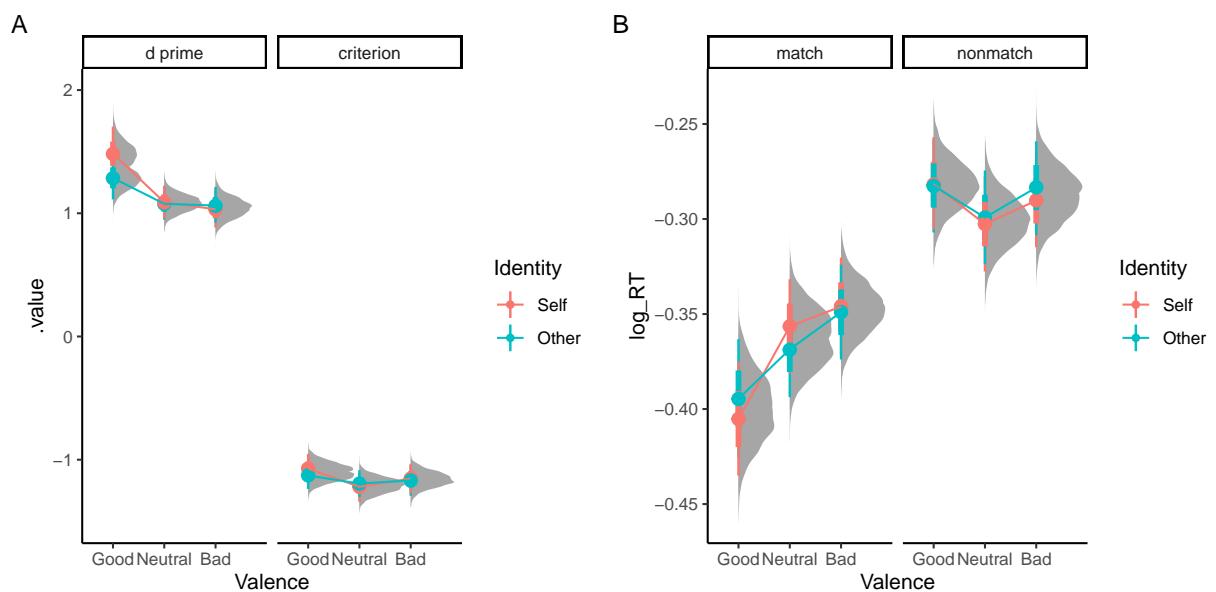


Figure 29. exp4b: Results of Bayesian GLM analysis.

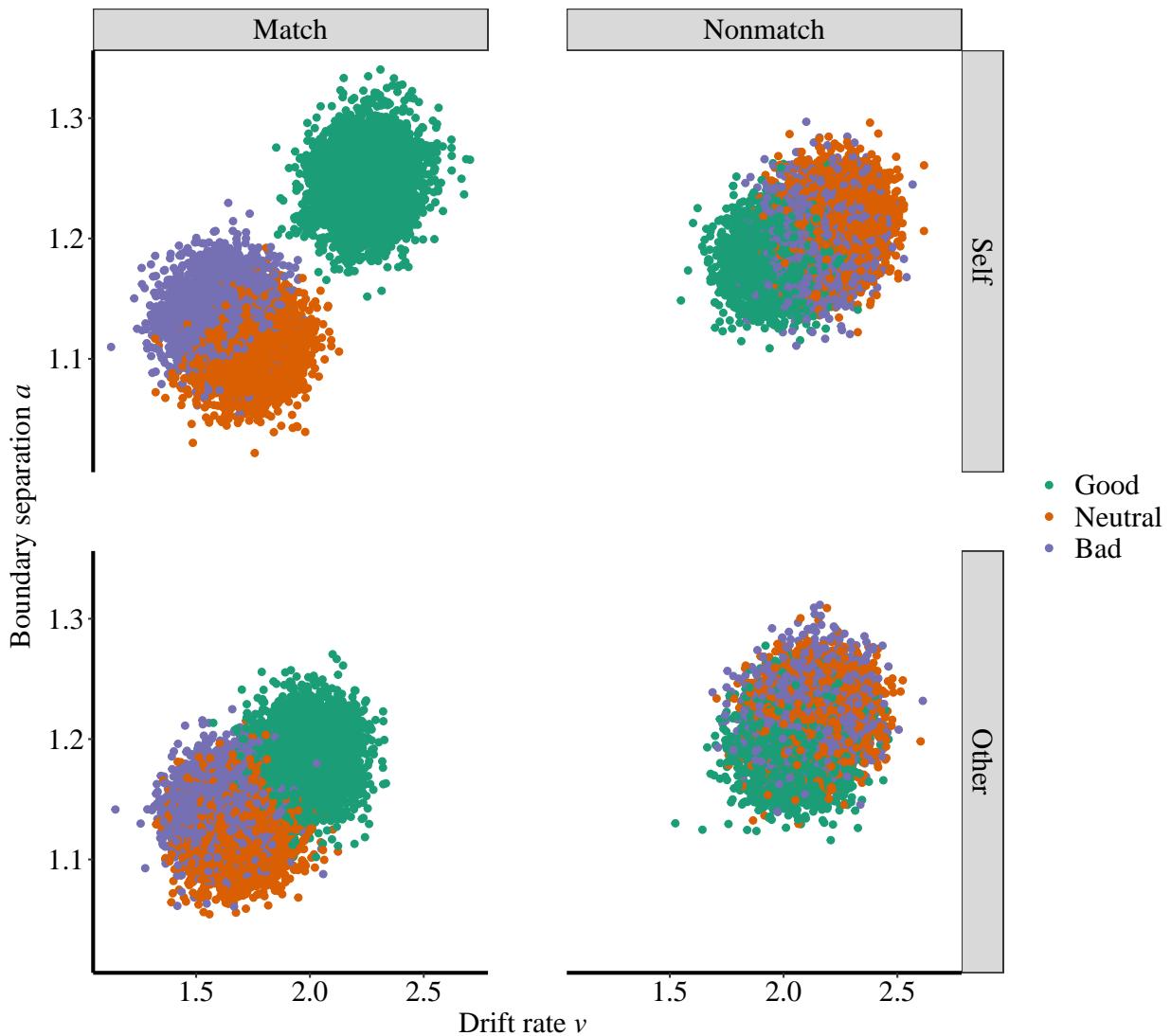


Figure 30. exp4b: Results of HDDM.

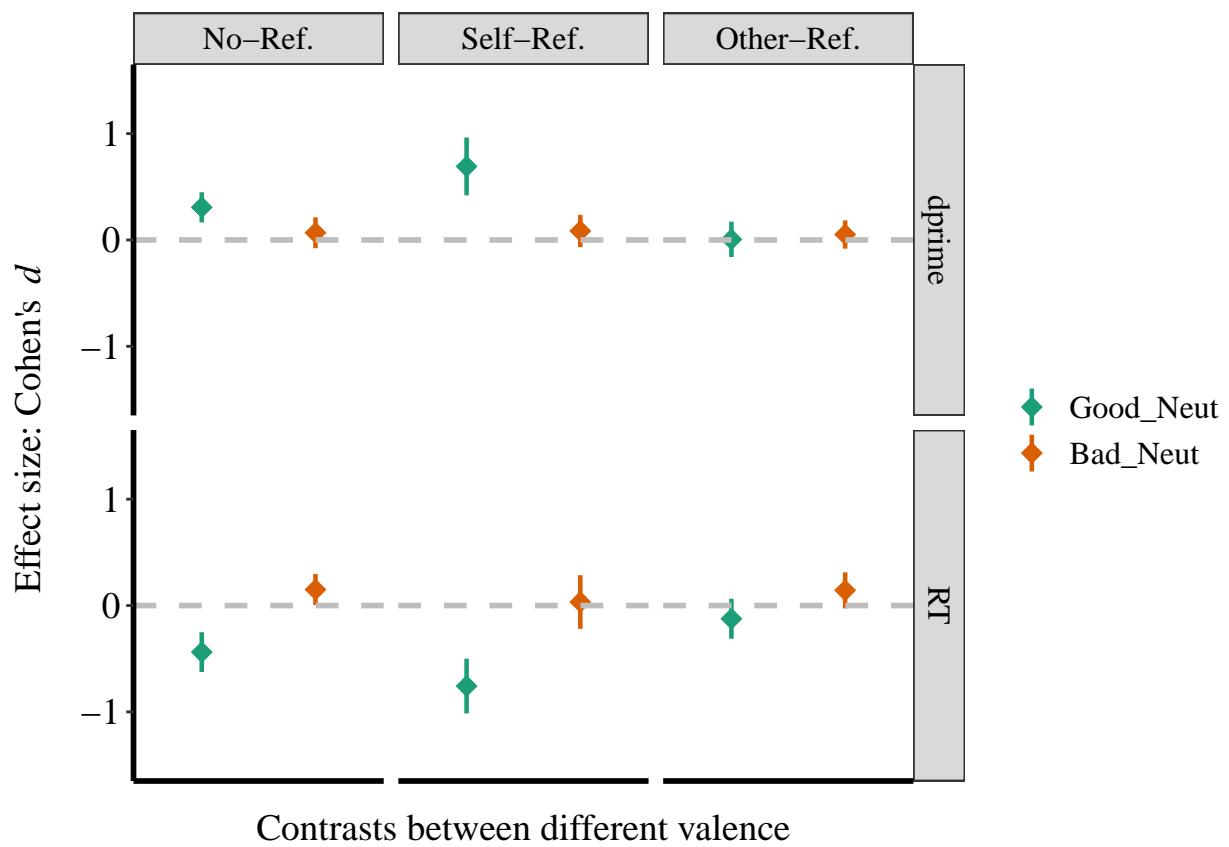


Figure 31. Effect size (Cohen's  $d$ ) of Valence.

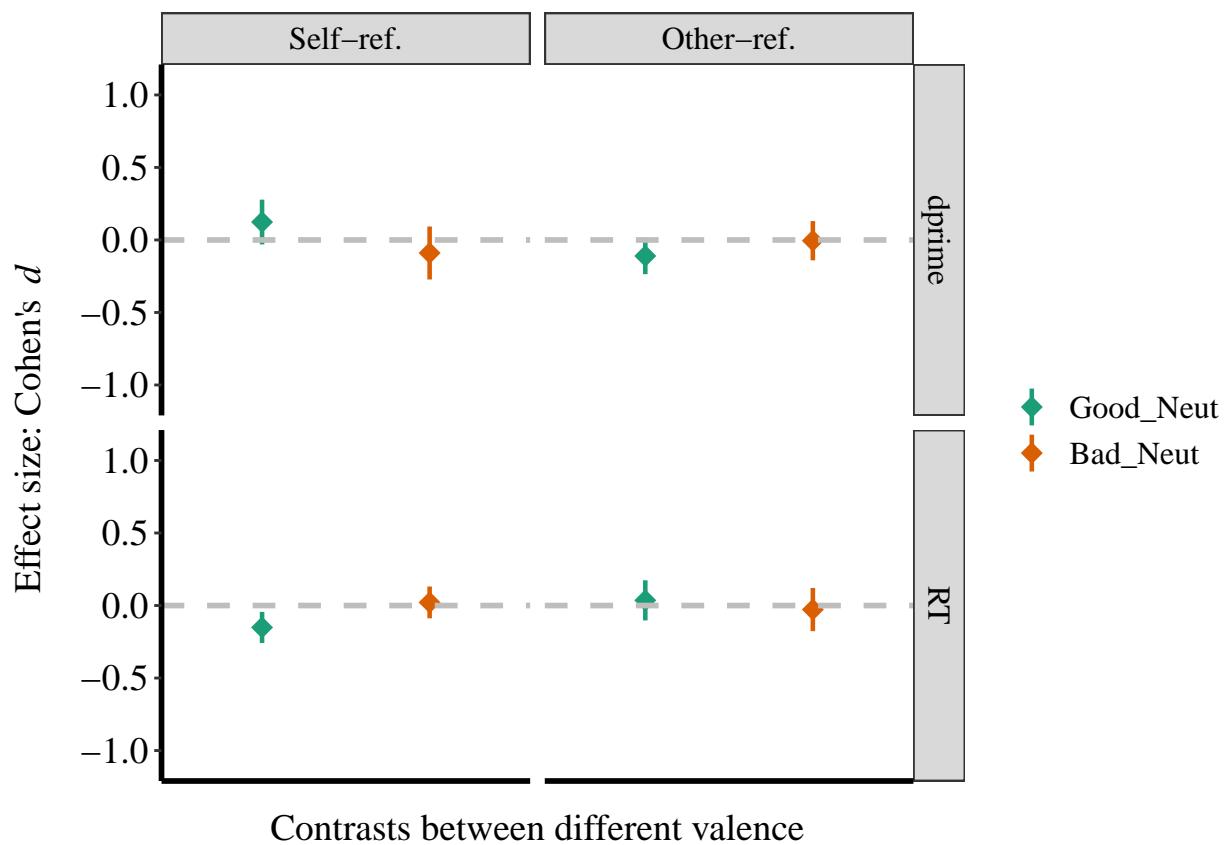


Figure 32. Effect size (Cohen's  $d$ ) of Valence in Exp4a.

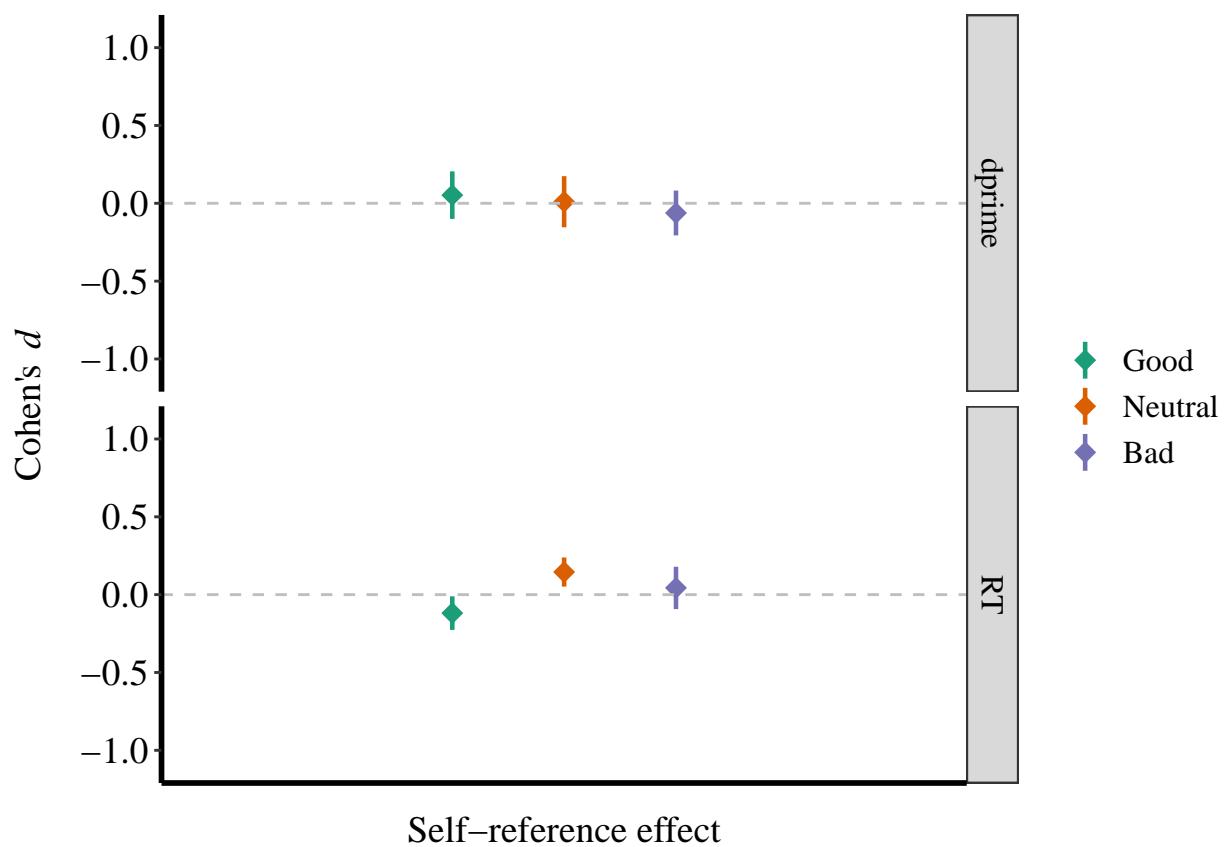


Figure 33. Effect size (Cohen's  $d$ ) of Valence in Exp4b.

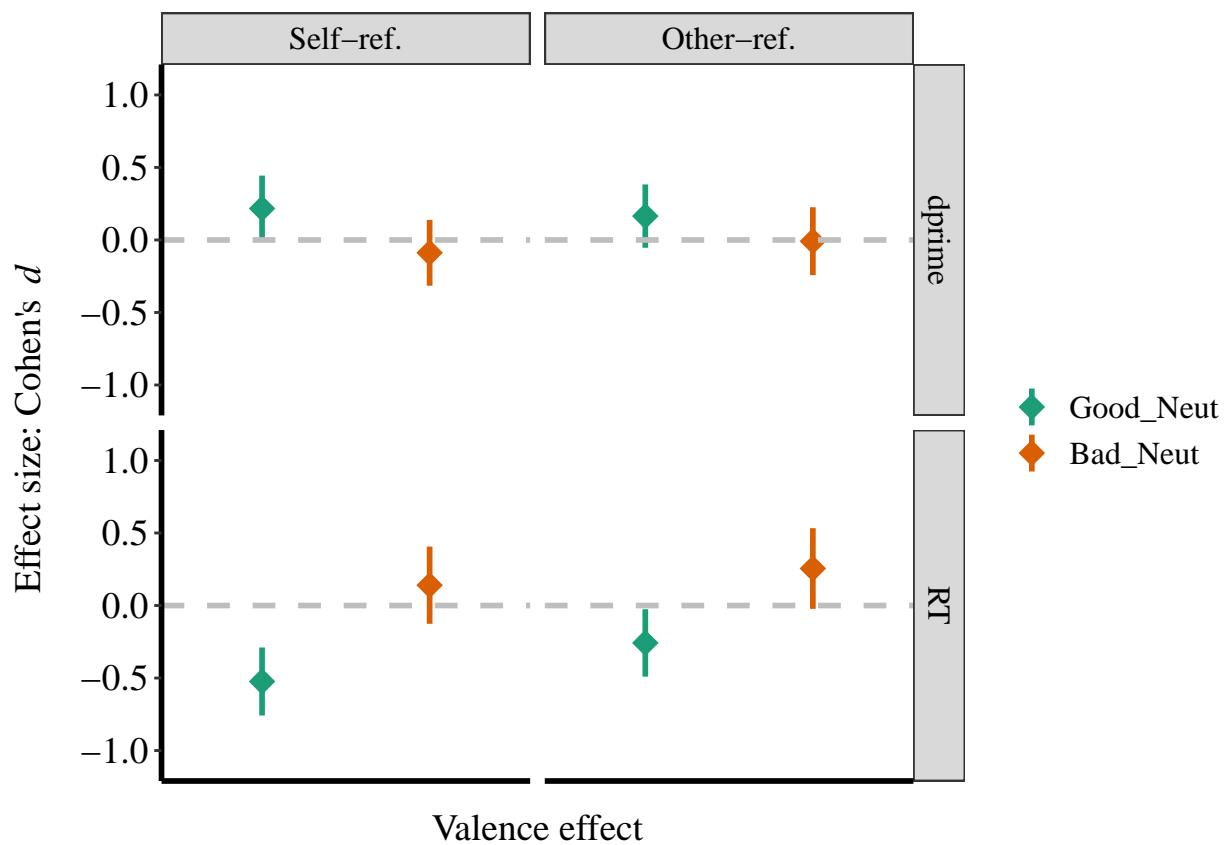


Figure 34. Effect size (Cohen's  $d$ ) of Valence in Exp4b.

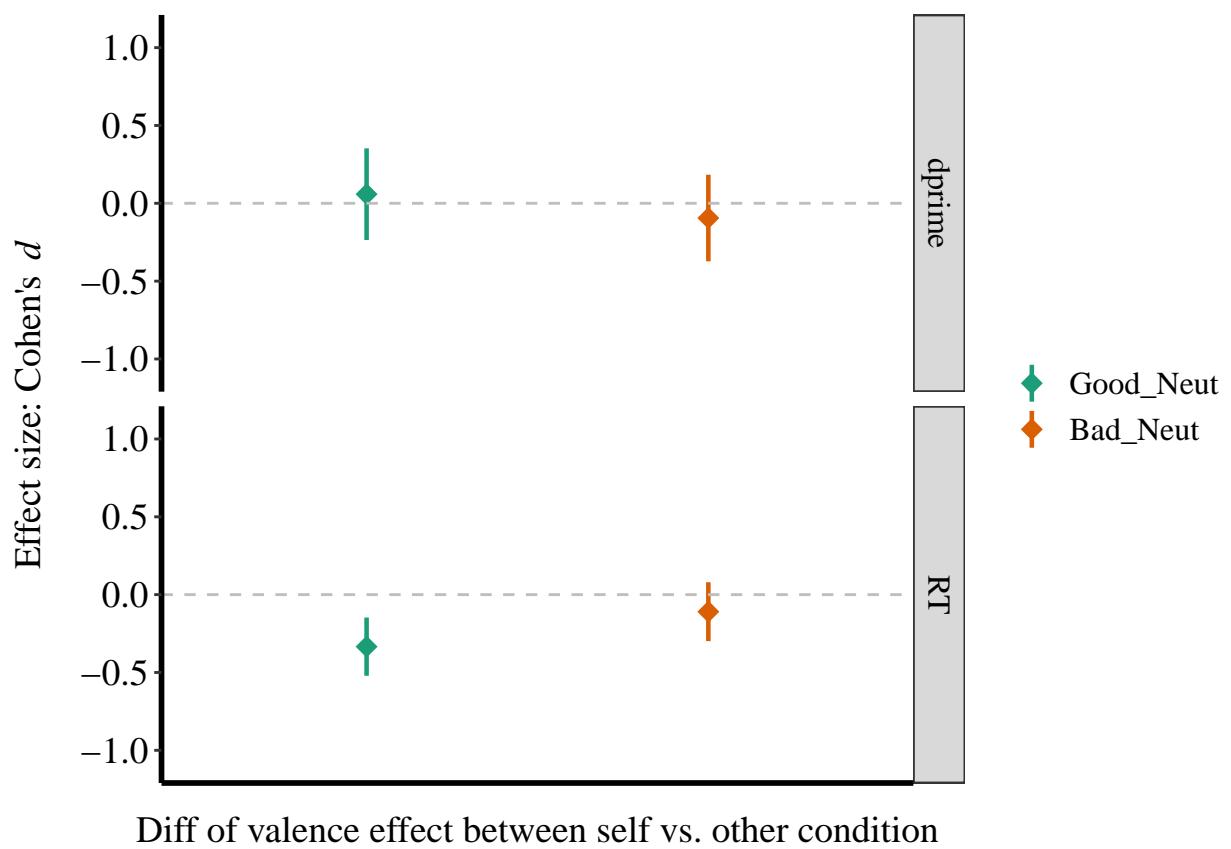


Figure 35. Effect size (Cohen's  $d$ ) of Valence in Exp4b.

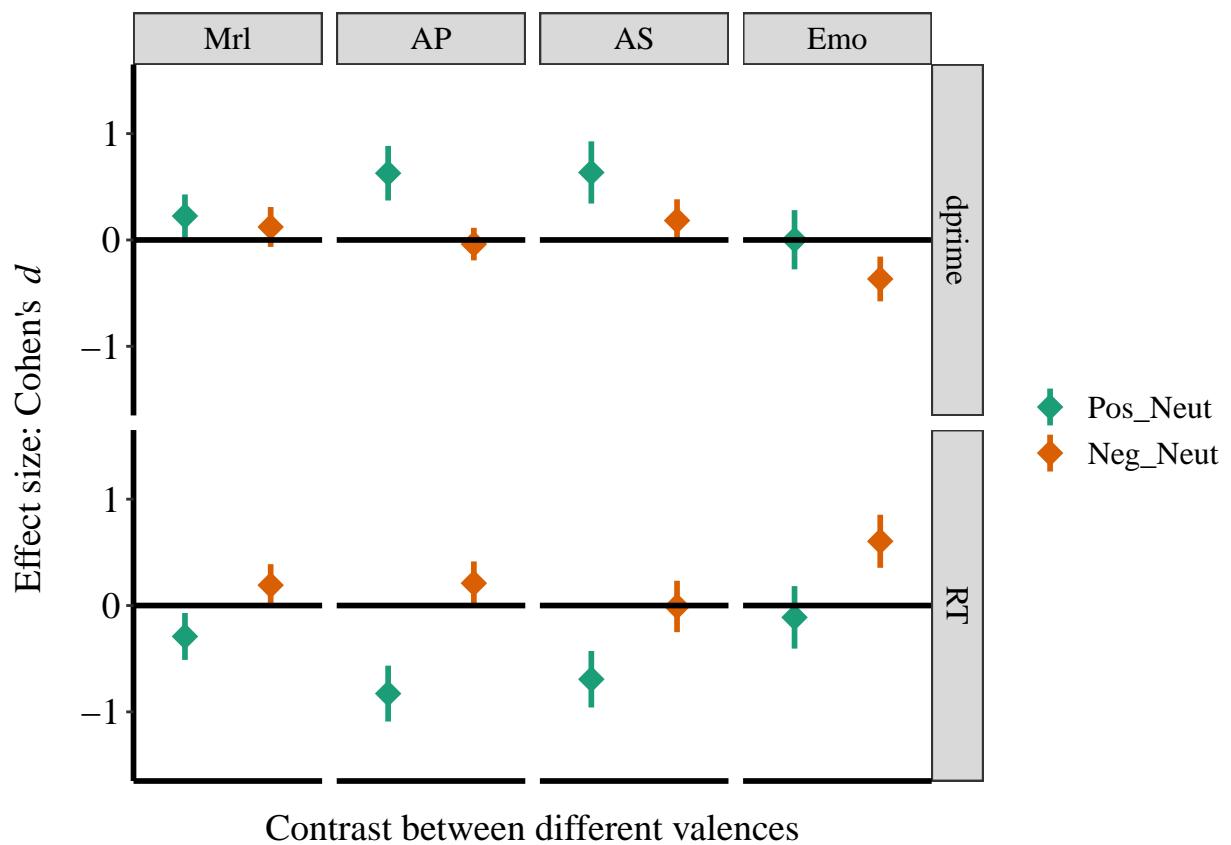


Figure 36. Effect size (Cohen's  $d$ ) of Valence in Exp5.

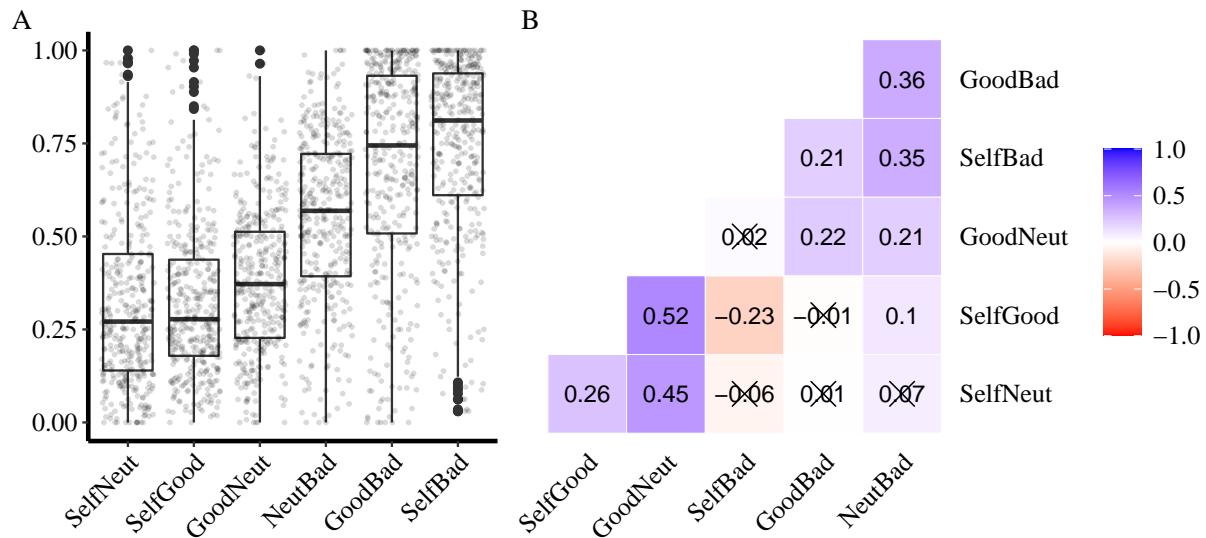


Figure 37. Self-rated personal distance

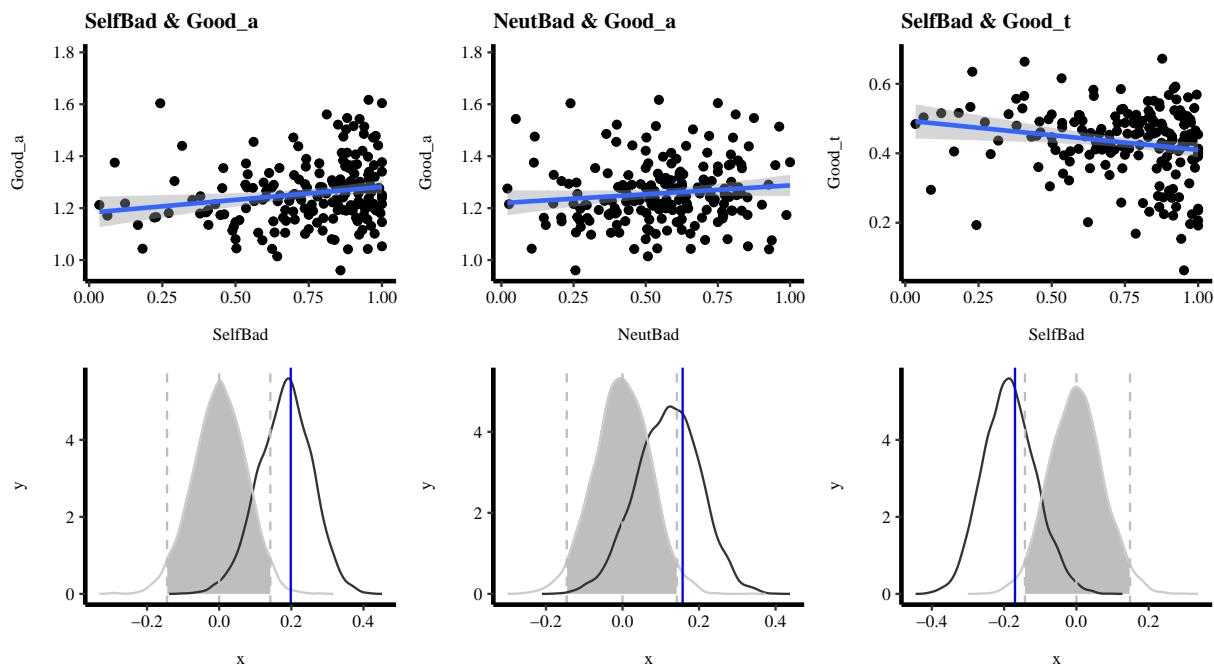
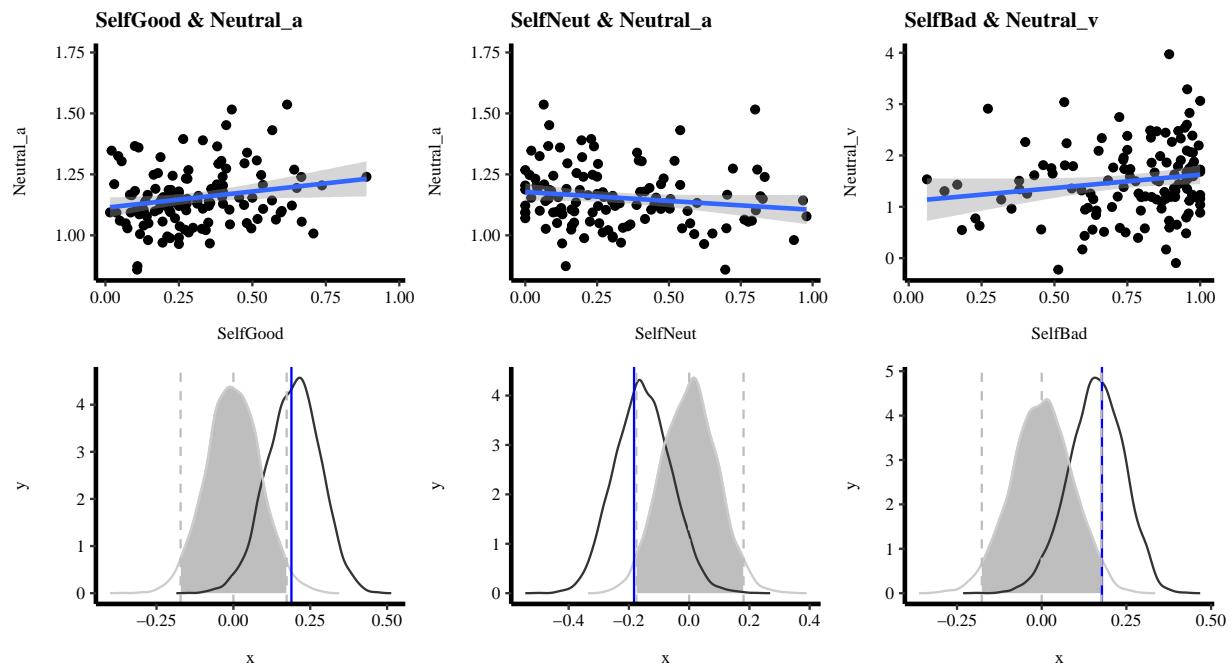


Figure 38. Correlation between moral identity and boundary separation of good condition; moral self-image and drift rate of good condition



*Figure 39.* Correlation between personal distance and boundary separation of neutral condition