

<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 [sentences in bracket are key ideas]

42 [Morality is the central of human social life]. People experience a substantial amount  
43 of moral events in everyday life (e.g., Hofmann, Wisneski, Brandt, & Skitka, 2014). When  
44 experiencing these events, it always involves judging “right” or “wrong”, “good” or “bad”.  
45 By judging “right” or “wrong”, people may implicitly infer “good” or “bad”, i.e., moral  
46 character (Uhlmann, Pizarro, & Diermeier, 2015). Similarly, moral character is a basic  
47 dimension of person perception (Abele, Ellemers, Fiske, Koch, & Yzerbyt, 2020; Goodwin,  
48 2015; Goodwin, Piazza, & Rozin, 2014; Willis & Todorov, 2006) and the most important  
49 aspect to evaluate the continuity of identity (Strohminger, Knobe, & Newman, 2017).

50 Given the importance of moral character, to successfully navigate in a social world, a  
51 person needs to both accurately evaluate others’ moral character and behave in a way that  
52 she/he is perceived as a moral person, or at least not a morally bad person. Maintaining a  
53 moral self-views is as important as making judgment about others’ moral character  
54 (Ellemers, Toorn, Paunov, & Leeuwen, 2019). Moral character is studied extensively both  
55 in person perception (Abele et al., 2020; Goodwin, 2015; Goodwin et al., 2014; Willis &  
56 Todorov, 2006) and moral self-view (Klein & Epley, 2016; Monin & Jordan, 2009;  
57 Strohminger et al., 2017; Tappin & McKay, 2017). Recent theorists are trying to bring  
58 them together and emphasize a person-centered moral psychology(Uhlmann et al., 2015).  
59 In this new perspective, role of perceiver’s self-relevance in morality has also been studied  
60 (e.g., Waytz, Dungan, & Young, 2013).

61 To date, however, as Freeman and Ambady (2011) put it, studies in the perception of  
62 moral character didn’t try to explain the perceptual process, rather, they are trying to  
63 explain the higher-order social cognitive processes that come after. Essentially, these

64 studies are perception of moral character without perceptual process. Without knowledge  
65 of perceptual processes, we can not have a full picture of how moral character is processed  
66 in our cognition. As an increasing attention is paid to perceptual process underlying social  
67 cognition, it's clear that perceptual processes are strongly influenced by social factors, such  
68 as group-categorization, stereotype (see Xiao, Coppin, & Bavel, 2016; Stolier & Freeman,  
69 2016). Given the importance of moral character and that moral character related  
70 information has strong influence on learning and memory (Carlson, Maréchal, Oud, Fehr,  
71 & Crockett, 2020; Stanley & De Brigard, 2019), one might expect that moral character  
72 related information could also play a role in perceptual process.

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

### 83 Perceptual process of moral character

84 [exp1a, b, c, and exp2]

85 [using associative learning task to study the moral character's influence on  
86 perception] Though it is theoretically possible that moral character related information  
87 may be prioritized in perceptual process, no empirical studies had directly explored this  
88 possibility. There were only a few studies about the temporal dynamics of judging the

<sup>89</sup> trustworthiness of face (e.g., Dzhelyova, Perrett, & Jentzsch, 2012), but trustworthy is not  
<sup>90</sup> equal to morality.

<sup>91</sup> One difficulty of studying the perceptual process of moral character is that moral  
<sup>92</sup> character is an inferred trait instead of observable feature. usually, one needs necessary  
<sup>93</sup> more sensory input, e.g., behavior history, to infer moral character of a person. For  
<sup>94</sup> example, Anderson et al. (2011) asked participant to first study the behavioral description  
<sup>95</sup> of faces and then asked them to perform a perceptual detection task. They assumed that  
<sup>96</sup> by learning the behavioral description of a person (represented by a face), participants can  
<sup>97</sup> acquire the moral related information about faces, and the associations could then bias the  
<sup>98</sup> perceptual processing of the faces (but see Stein, Grubb, Bertrand, Suh, and Verosky  
<sup>99</sup> (2017)). One drawback of this approach is that participants may differ greatly when  
<sup>100</sup> inferring the moral character of the person from behavioral descriptions, given that notion  
<sup>101</sup> what is morality itself is varying across population (Henrich, Heine, & Norenzayan, 2010)  
<sup>102</sup> and those descriptions and faces may themselves are idiosyncratic, therefore, introduced  
<sup>103</sup> large variation in experimental design.

<sup>104</sup> An alternative is to use abstract semantic concepts. Abstract concepts of moral  
<sup>105</sup> character are used to describe and represent moral characters. These abstract concepts  
<sup>106</sup> may be part of a dynamic network in which sensory cue, concrete behaviors and other  
<sup>107</sup> information can activate/inhibit each other (e.g., aggressiveness) (Amodio, 2019; Freeman  
<sup>108</sup> & Ambady, 2011). If a concept of moral character (e.g., good person) is activated, it  
<sup>109</sup> should be able to influence on the perceptual process of the visual cues through the  
<sup>110</sup> dynamic network, especially when the perceptual decision-making is about the concept-cue  
<sup>111</sup> association. In this case, abstract concepts of moral character may serve as signal of moral  
<sup>112</sup> reputation (for others) or moral self-concept. Indeed, previous studies used the moral  
<sup>113</sup> words and found that moral related information can be perceived faster (Gantman & Van  
<sup>114</sup> Bavel, 2014, but see, @firestone\_enhanced\_2015). If moral character is an important in  
<sup>115</sup> person perception, then, just as those other information such as races and stereotype (see

<sup>116</sup> Xiao et al., 2016), moral character related concept might change the perceptual processes.

<sup>117</sup> To investigate the above possibility, we used an associative learning paradigm to  
<sup>118</sup> study how moral character concept change perceptual decision-making. In this paradigm,  
<sup>119</sup> simple geometric shapes were paired with different words whose dominant meaning is  
<sup>120</sup> describing the moral character of a person. Participants first learn the associations between  
<sup>121</sup> shapes and words, e.g., triangle is a good-person. After building direct association between  
<sup>122</sup> the abstract moral characters and visual cues, participants then perform a matching task  
<sup>123</sup> to judge whether the shape-word pair presented on the screen match the association they  
<sup>124</sup> learned. This paradigm has been used in studying the perceptual process of self-concept,  
<sup>125</sup> but had also proven useful in studying other concepts like social group (Enock, Hewstone,  
<sup>126</sup> Lockwood, & Sui, 2020; Enock, Sui, Hewstone, & Humphreys, 2018). By using simple and  
<sup>127</sup> morally neutral shapes, we controlled the variations caused by visual cues.

<sup>128</sup> Our first question is, whether the words used the in the associative paradigm is really  
<sup>129</sup> related to the moral character? As we reviewed above, previous theories, especially the  
<sup>130</sup> interactive dynamic theory, would support this assumption. To validate that moral  
<sup>131</sup> character concepts activated moral character as a social cue, we used four experiments to  
<sup>132</sup> explore and validate the paradigm. The first experiment directly adopted associative  
<sup>133</sup> paradigm and change the words from “self”, “friend”, and “stranger” to “good-person”,  
<sup>134</sup> “neutral-person”, and “bad-person”. Then, we change the words to the ones that have  
<sup>135</sup> more explicit moral meaning (“kind-person”, “neutral-person”, and “evil-person”). Then,  
<sup>136</sup> as in Anderson et al. (2011), we asked participant to learn the association between three  
<sup>137</sup> different behavioral histories and three different names, and then use the names, as moral  
<sup>138</sup> character words, for associative learning. Finally, we also tested that simultaneously  
<sup>139</sup> present shape-word pair and sequentially present word and shape didn’t change the  
<sup>140</sup> pattern. All of these four experiments showed a robust effect of moral character, that is,  
<sup>141</sup> the positive moral character associated stimuli were prioritized.

<sup>142</sup> **Morality as a social-categorization?**

<sup>143</sup> [possible explanations: person-based self-categorization vs. stimuli-based valence] The  
<sup>144</sup> robust pattern from our first four experiment suggested that there are some reliable  
<sup>145</sup> mechanisms underneath the effect. One possible explanation is the value-based attention,  
<sup>146</sup> which suggested that valuable stimuli is prioritized in our low-level cognitive processes.  
<sup>147</sup> Because positive moral character is potentially rewarding, e.g., potential cooperators, it is  
<sup>148</sup> valuable to individuals and therefore being prioritized. There are also evidence consistent  
<sup>149</sup> with this idea []. For example, XXX found that trustworthy faces attracted attention more  
<sup>150</sup> than untrustworthy faces, probably because trustworthy faces are more likely to be the  
<sup>151</sup> collaborative partners subsequent tasks, which will bring reward. This explanation has an  
<sup>152</sup> implicit assumption, that is, participants were automatically viewing these stimuli as  
<sup>153</sup> self-relevant (Juechems & Summerfield, 2019; Reicher & Hopkins, 2016) and  
<sup>154</sup> threatening/rewarding because of their semantic meaning. In this explanation, we will view  
<sup>155</sup> the moral concept, and the moral character represented by the concept, as objects and only  
<sup>156</sup> judge whether they are rewarding/threatening or potentially rewarding/threatening to us.

<sup>157</sup> Another possibility is that we will perceive those moral character as person and  
<sup>158</sup> automatic categorize whether they are ingroup or ougroup, that is, the social  
<sup>159</sup> categorization process. This account assumed that moral character served as a way to  
<sup>160</sup> categorize other. In the first four experiments' situation, the identity of the moral  
<sup>161</sup> character is ambiguous, participants may automatically categorize morally good people as  
<sup>162</sup> ingroup and therefore preferentially processed these information.

<sup>163</sup> However, the above four experiments can not distinguish between these two  
<sup>164</sup> possibilities, because the concept “good-peron” can both be rewarding and be categorized  
<sup>165</sup> as ingroup memeber, and previous studies using associative learning paradigm revealed  
<sup>166</sup> that both rewarding stimuli (e.g., Sui, He, & Humphreys, 2012) and in-group information  
<sup>167</sup> [Enock et al. (2018); enock\_overlap\_2020] are prioritized.

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

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

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

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

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

226 Key concepts and discussing points:

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

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

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

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

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

237 Morality as a way for social-categorization (McHugh, McGann, Igou, & Kinsella,  
238 2019)? People are more likely to identify themselves with trustworthy faces (Verosky &  
239 Todorov, 2010) (trustworthy faces has longer RTs).

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

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

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

246

## Disclosures

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

251 All the experiments reported were not pre-registered. Most experiments (1a ~ 6b,  
252 except experiment 3b) reported in the current study were first finished between 2014 to  
253 2016 in Tsinghua University, Beijing, China. Participants in these experiments were  
254 recruited in the local community. To increase the sample size of experiments to 50 or more  
255 (Simmons, Nelson, & Simonsohn, 2013), we recruited additional participants in Wenzhou  
256 University, Wenzhou, China in 2017 for experiment 1a, 1b, 4a, and 4b. Experiment 3b was  
257 finished in Wenzhou University in 2017. To have a better estimation of the effect size, we  
258 included the data from two experiments (experiment 7a, 7b) that were reported in Hu et  
259 al. (2020) (See Table S1 for overview of these experiments).

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

262

## General methods

### 263 Design and Procedure

264 This series of experiments started to test the effect of instantly acquired true self  
265 (moral self) on perceptual decision-making. For this purpose, we used the social associative  
266 learning paradigm (or tagging paradigm)(Sui et al., 2012), in which participants first  
267 learned the associations between geometric shapes and labels of person with different moral  
268 character (e.g., in first three studies, the triangle, square, and circle and good person,  
269 neutral person, and bad person, respectively). The associations of the shapes and label

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

276 The experiment 1a, 1b, 1c, 2, and 6a shared a 2 (matching: match vs. nonmatch) by  
277 3 (moral valence: good vs. neutral vs. bad) within-subject design. Experiment 1a was the  
278 first one of the whole series studies and 1b, 1c, and 2 were conducted to exclude the  
279 potential confounding factors. More specifically, experiment 1b used different Chinese  
280 words as label to test whether the effect only occurred with certain familiar words.  
281 Experiment 1c manipulated the moral valence indirectly: participants first learned to  
282 associate different moral behaviors with different neutral names, after remembered the  
283 association, they then performed the perceptual matching task by associating names with  
284 different shapes. Experiment 2 further tested whether the way we presented the stimuli  
285 influence the effect of valence, by sequentially presenting labels and shapes. Note that part  
286 of participants of experiment 2 were from experiment 1a because we originally planned a  
287 cross task comparison. Experiment 6a, which shared the same design as experiment 2, was  
288 an EEG experiment which aimed at exploring the neural correlates of the effect. But we  
289 will focus on the behavioral results of experiment 6a in the current manuscript.

290 For experiment 3a, 3b, 4a, 4b, 6b, 7a, and 7b, we included self-reference as another  
291 within-subject variable in the experimental design. For example, the experiment 3a directly  
292 extend the design of experiment 1a into a 2 (matchness: match vs. nonmatch) by 2  
293 (reference: self vs. other) by 3 (moral valence: good vs. neutral vs. bad) within-subject  
294 design. Thus in experiment 3a, there were six conditions (good-self, neutral-self, bad-self,  
295 good-other, neutral-other, and bad-other) and six shapes (triangle, square, circle, diamond,  
296 pentagon, and trapezoids). The experiment 6b was an EEG experiment extended from

297 experiment 3a but presented the label and shape sequentially. Because of the relatively  
298 high working memory load (six label-shape pairs), experiment 6b were conducted in two  
299 days: the first day participants finished perceptual matching task as a practice, and the  
300 second day, they finished the task again while the EEG signals were recorded. Experiment  
301 3b was designed to separate the self-referential trials and other-referential trials. That is,  
302 participants finished two different blocks: in the self-referential blocks, they only responded  
303 to good-self, neutral-self, and bad-self, with half match trials and half non-match trials; for  
304 the other-reference blocks, they only responded to good-other, neutral-other, and  
305 bad-other. Experiment 7a and 7b were designed to test the cross task robustness of the  
306 effect we observed in the aforementioned experiments (see, Hu et al., 2020). The matching  
307 task in these two experiments shared the same design with experiment 3a, but only with  
308 two moral valence, i.e., good vs. bad. We didn't include the neutral condition in  
309 experiment 7a and 7b because we found that the neutral and bad conditions constantly  
310 showed non-significant results in experiment 1 ~ 6.

311 Experiment 4a and 4b were design to test the automaticity of the binding between  
312 self/other and moral valence. In 4a, we used only two labels (self vs. other) and two shapes  
313 (circle, square). To manipulate the moral valence, we added the moral-related words within  
314 the shape and instructed participants to ignore the words in the shape during the task. In  
315 4b, we reversed the role of self-reference and valence in the task: participant learnt three  
316 labels (good-person, neutral-person, and bad-person) and three shapes (circle, square, and  
317 triangle), and the words related to identity, “self” or “other”, were presented in the shapes.  
318 As in 4a, participants were told to ignore the words inside the shape during the task.

319 Finally, experiment 5 was design to test the specificity of the moral valence. We  
320 extended experiment 1a with an additional independent variable: domains of the valence  
321 words. More specifically, besides the moral valence, we also added valence from other  
322 domains: appearance of person (beautiful, neutral, ugly), appearance of a scene (beautiful,  
323 neutral, ugly), and emotion (happy, neutral, and sad). Label-shape pairs from different

324 domains were separated into different blocks.

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

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

#### 344 Data analysis

345 **Analysis of individual study.** We used the `tidyverse` of r (see script  
346 `Load_save_data.r`) to exclude the practicing trials, invalid trials of each participants, and  
347 invalid participants, if there were any, in the raw data. Results of each experiment were  
348 then analyzed in three different approaches.

349        ***Classic NHST.***

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

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

364        where the  $HR$  means hit rate and the  $FAR$  mean false alarm rate.  $zHR$  and  $zFAR$  are  
 365        the standardized hit rate and false alarm rates, respectively. These two  $z$ -scores were  
 366        converted from proportion (i.e., hit rate or false alarm rate) by inverse cumulative normal  
 367        density function,  $\Phi^{-1}$  ( $\Phi$  is the cumulative normal density function, and is used convert  $z$   
 368        score into probabilities). Another parameter of signal detection theory, response criterion  $c$ ,  
 369        is defined by the negative standardized false alarm rate (DeCarlo, 1998):  $-zFAR$ .

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

374 To control the false positive rate when conducting the post-hoc comparisons, we used

375 Bonferroni correction.

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

377 The classic NHST approach may ignore the uncertainty in estimate of the parameters

378 for SDT (Rouder & Lu, 2005), and using mean RT assumes normal distribution of RT

379 data, which is always not true because RTs distribution is skewed (Rousselet & Wilcox,

380 2019). To better estimate the uncertainty and use a more appropriate model, we also tried

381 Bayesian hierarchical generalized linear model to analyze each experiment's accuracy and

382 RTs data. We used BRMs (Bürkner, 2017) to build the model, which used Stan (Carpenter

383 et al., 2017) to estimate the posterior.

384 In the GLM model, we assume that the accuracy of each trial is Bernoulli distributed

385 (binomial with 1 trial), with probability  $p_i$  that  $y_i = 1$ .

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

386 In the perceptual matching task, the probability  $p_i$  can then be modeled as a function of

387 the trial type:

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

388 The outcomes  $y_i$  are 0 if the participant responded "nonmatch" on trial  $i$ , 1 if they

389 responded "match". The probability of the "match" response for trial  $i$  for a participant is

390  $p_i$ . We then write the generalized linear model on the probits (z-scores;  $\Phi$ , "Phi") of  $ps$ .  $\Phi$

391 is the cumulative normal density function and maps  $z$  scores to probabilities. Given this

392 parameterization, the intercept of the model ( $\beta_0$ ) is the standardized false alarm rate

393 (probability of saying 1 when predictor is 0), which we take as our criterion  $c$ . The slope of

394 the model ( $\beta_1$ ) is the increase of saying 1 when predictor is 1, in  $z$ -scores, which is another

395 expression of  $d'$ . Therefore,  $c = -z\text{HR} = -\beta_0$ , and  $d' = \beta_1$ .

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

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

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

$$\Phi(p_{ij}) = \beta_{0j} + \beta_{1j} \text{IsMatch}_{ij} * \text{Valence}_{ij}$$

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

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

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

407 For the reaction time, we used the log normal distribution  
 408 ([https://lindeloev.github.io/shiny-rt/#34\\_\(shifted\)\\_log-normal](https://lindeloev.github.io/shiny-rt/#34_(shifted)_log-normal)). This distribution has  
 409 two parameters:  $\mu$ ,  $\sigma$ .  $\mu$  is the mean of the logNormal distribution, and  $\sigma$  is the disperse of  
 410 the distribution. The log normal distribution can be extended to shifted log normal  
 411 distribution, with one more parameter: shift, which is the earliest possible response.

$$y_i = \beta_0 + \beta_1 * \text{IsMatch}_i * \text{Valence}_i$$

412 Shifted log-normal distribution:

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

<sup>413</sup>  $y_{ij}$  is the RT of the  $i$ th trial of the  $j$ th participants.

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

$$\sigma_j \sim Cauchy()$$

<sup>414</sup> **Hierarchical drift diffusion model (HDDM).**

<sup>415</sup> To further explore the psychological mechanism under perceptual decision-making, we  
<sup>416</sup> used HDDM (Wiecki, Sofer, & Frank, 2013) to model our RTs and accuracy data. We used  
<sup>417</sup> the prior implemented in HDDM, that is, informative priors that constrains parameter  
<sup>418</sup> estimates to be in the range of plausible values based on past literature (Matzke &  
<sup>419</sup> Wagenmakers, 2009). As reported in Hu et al. (2020), we used the response code approach,  
<sup>420</sup> match response were coded as 1 and nonmatch responses were coded as 0. To fully explore  
<sup>421</sup> all parameters, we allow all four parameters of DDM free to vary. We then extracted the  
<sup>422</sup> estimation of all the four parameters for each participants for the correlation analyses.  
<sup>423</sup> However, because the starting point is only related to response (match vs. non-match) but  
<sup>424</sup> not the valence of the stimuli, we didn't included it in correlation analysis.

<sup>425</sup> **Synthesized results.** We also reported the synthesized results from the  
<sup>426</sup> experiments, because many of them shared the similar experimental design. We reported  
<sup>427</sup> the results in five parts: valence effect, explicit interaction between valence and  
<sup>428</sup> self-relevance, implicit interaction between valence and self-relevance, specificity of valence  
<sup>429</sup> effect, and behavior-questionnaire correlation.

<sup>430</sup> For the first two parts, we reported the synthesized results from Frequentist's  
<sup>431</sup> approach(mini-meta-analysis, Goh, Hall, & Rosenthal, 2016). The mini meta-analyses were  
<sup>432</sup> carried out by using `metafor` package (Viechtbauer, 2010). We first calculated the mean of  
<sup>433</sup>  $d'$  and RT of each condition for each participant, then calculate the effect size (Cohen's  $d$ )

434 and variance of the effect size for all contrast we interested: Good v. Bad, Good v.  
 435 Neutral, and Bad v. Neutral for the effect of valence, and self vs. other for the effect of  
 436 self-relevance. Cohen's  $d$  and its variance were estimated using the following formula  
 437 (Cooper, Hedges, & Valentine, 2009):

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

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

438  $M_1$  is the mean of the first condition,  $sd_1$  is the standard deviation of the first  
 439 condition, while  $M_2$  is the mean of the second condition,  $sd_2$  is the standard deviation of  
 440 the second condition.  $r$  is the correlation coefficient between data from first and second  
 441 condition.  $n$  is the number of data point (in our case the number of participants included  
 442 in our research).

443 The effect size from each experiment were then synthesized by random effect model  
 444 using `metafor` (Viechtbauer, 2010). Note that to avoid the cases that some participants  
 445 participated more than one experiments, we inspected the all available information of  
 446 participants and only included participants' results from their first participation. As  
 447 mentioned above, 24 participants were intentionally recruited to participate both exp 1a  
 448 and exp 2, we only included their results from experiment 1a in the meta-analysis.

449 We also estimated the synthesized effect size using Bayesian hierarchical model,  
 450 which extended the two-level hierarchical model in each experiment into three-level model,  
 451 which experiment as an additional level. For SDT, we can use a nested hierarchical model  
 452 to model all the experiment with similar design:

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

<sup>453</sup> where

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

<sup>454</sup> The outcomes  $y_{ijk}$  are 0 if participant  $j$  in experiment  $k$  responded “nonmatch” on trial  $i$ ,

<sup>455</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>456</sup> and the experiment level parameter  $mu_{0k}$  and  $mu_{1k}$  is from a higher order

<sup>457</sup> distribution:

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

<sup>458</sup> in which  $\mu_0$  and  $\mu_1$  means the population level parameter.

<sup>459</sup> This model can be easily expand to three-level model in which participants and

<sup>460</sup> experiments are two group level variable and participants were nested in the experiments.

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

<sup>461</sup>  $y_{ijk}$  is the RT of the  $i$ th trial of the  $j$ th participants in the  $k$ th experiment.

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

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

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

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

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

466 ***Valence effect.***

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

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

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

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

474 In the third part, we focused on experiment 4a and 4b, which were designed to  
475 examine the implicit effect of the interaction between moral valence and self-referential  
476 processing. We are interested in one particular question: will self-referential and morally  
477 positive valence had a mutual facilitation effect. That is, when moral valence (experiment  
478 4a) or self-referential (experiment 4a) was presented as task-irrelevant stimuli, whether  
479 they would facilitate self-referential or valence effect on perceptual decision-making. For  
480 experiment 4a, we reported the comparisons between different valence conditions under the  
481 self-referential task and other-referential task. For experiment 4b, we first calculated the  
482 effect of valence for both self- and other-referential conditions and then compared the effect  
483 size of these three contrast from self-referential condition and from other-referential  
484 condition. Note that the results were also analyzed in a standard repeated measure  
485 ANOVA (see supplementary materials).

486        ***Specificity of the valence effect.***

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

491        ***Behavior-Questionnaire correlation.***

492        Finally, we explored correlation between results from behavioral results and  
493        self-reported measures.

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

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

504        Based on results form the experiment, we reason that the correlation between  
505        behavioral result in self-referential will appear in the data without mentioning the  
506        self/other (exp 3a, 3b, 6a, 7a, 7b). To this end, we first calculated the correlation between  
507        behavioral indicators and questionnaires for self-referential and other-referential separately.  
508        Given the small sample size of the data ( $N =$ ), we used a relative liberal threshold for  
509        these exploration ( $\alpha = 0.1$ ).

510        Then we confirmed the significant results from the data without self- and  
511        other-referential (exp1a, 1b, 1c, 2, 5, 6a). In this confirmatory analysis, we used  $\alpha =$

512 0.05 and used bootstrap by BootES package (Kirby & Gerlanc, 2013) to estimate the  
513 correlation. To avoid false positive, we further determined the threshold for significant by  
514 permutation. More specifically, for each pairs that initially with  $p < .05$ , we randomly  
515 shuffle the participants data of each score and calculated the correlation between the  
516 shuffled vectors. After repeating this procedure for 5000 times, we choose arrange these  
517 5000 correlation coefficients and use the 95% percentile number as our threshold.

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

519 In this part, we report five experiments that aimed at testing whether an associative  
520 learning task, in which concepts of moral character are associated with geometric shapes,  
521 will impact the perceptual decision-making.

522 **Experiment 1a**

523 **Methods.**

524 ***Participants.***

525 57 college students (38 female, age =  $20.75 \pm 2.54$  years) participated. 39 of them  
526 were recruited from Tsinghua University community in 2014; 18 were recruited from  
527 Wenzhou University in 2017. All participants were right-handed except one, and all had  
528 normal or corrected-to-normal vision. Informed consent was obtained from all participants  
529 prior to the experiment according to procedures approved by the local ethics committees. 6  
530 participant's data were excluded from analysis because nearly random level of accuracy,  
531 leaving 51 participants (34 female, age =  $20.72 \pm 2.44$  years).

532 ***Stimuli and Tasks.***

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

536      ***Procedure.***

537      This experiment had two phases. First, there was a brief learning stage. Participants  
538      were asked to learn the relationship between geometric shapes (triangle, square, and circle)  
539      and different concepts of moral character (bad person, a good person, or a neutral person).  
540      For example, a participant was told, “bad person is a circle; good person is a triangle; and  
541      a neutral person is a square.” After participant remembered the associations (usually in a  
542      few minutes), participants started a practicing phase of matching task which has the exact  
543      task as in the experimental task.

544      In the experimental task, participants judged whether shape–label pairs, which were  
545      subsequently presented, were correct (i.e., the same as they learned). Each trial started  
546      with the presentation of a central fixation cross for 500 ms. Subsequently, a pairing of a  
547      shape and label (good person, bad person, and neutral person) was presented for 100 ms.  
548      The pair presented could confirm to the verbal instruction for each pairing given in the  
549      training stage, or it could be a recombination of a shape with a different label, with the  
550      shape–label pairings being generated at random. The next frame showed a blank for  
551      1100ms. Participants were expected to judge whether the shape was correctly assigned to  
552      the person by pressing one of the two response buttons as quickly and accurately as  
553      possible within this timeframe (to encourage immediate responding). Feedback (correct or  
554      incorrect) was given on the screen for 500 ms at the end of each trial, if no response  
555      detected, “too slow” was presented to remind participants to accelerate. Participants were  
556      informed of their overall accuracy at the end of each block. The practice phase finished and  
557      the experimental task began after the overall performance of accuracy during practice  
558      phase achieved 60%.

559      For participants from the Tsinghua community, they completed 6 experimental blocks  
560      of 60 trials. Thus, there were 60 trials in each condition (bad-person match, bad-person  
561      nonmatch, good-person match, good-person nonmatch, neutral-person match, and

562 neutral-person nonmatch). For the participants from Wenzhou University, they finished 6  
563 blocks of 120 trials, therefore, 120 trials for each condition.

564 ***Data analysis.***

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

568 **Results.**

569 ***Classic NHST.***

570 *d prime.*

571 Figure ?? shows *d'* and reaction times during the perceptual matching task. We  
572 conducted a single factor (valence: good, neutral, bad) repeated measure ANOVA.

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

580 *Reaction times.*

581 We conducted 2 (Matchness: match v. nonmatch) by 3 (Valence: good, neutral, bad)  
582 repeated measure ANOVA. We found the main effect of Matchness ( $F(1, 50) = 232.39$ ,  
583  $MSE = 948.92$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .104$ ), main effect of valence ( $F(1.87, 93.31) = 9.62$ ,  
584  $MSE = 1,673.86$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .016$ ), and interaction between Matchness and Valence  
585 ( $F(1.73, 86.65) = 8.52$ ,  $MSE = 1,441.75$ ,  $p = .001$ ,  $\hat{\eta}_G^2 = .011$ ).

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

587 matched trials, we found the effect of valence . We further examined the effect of valence  
588 for both self and other for matched trials. We found that shapes associated with Good  
589 Person (684 ms, SE = 11.5) responded faster than Neutral (709 ms, SE = 11.5),  $t(50) =$   
590 -2.265,  $p = 0.0702$ ) and Bad Person (728 ms, SE = 11.7),  $t(50) = -4.41$ ,  $p = 0.0002$ ), and  
591 the Neutral condition was faster than the Bad condition,  $t(50) = -2.495$ ,  $p = 0.0415$ ). For  
592 non-matched trials, there was no significant effect of Valence ()�.

593 ***Bayesian hierarchical GLM.***

594 *d prime.*

595 We fitted a Bayesian hierarchical GLM for signal detection theory. The results  
596 showed that when the shapes were tagged with labels with different moral valence, the  
597 sensitivity ( $d'$ ) and criteria ( $c$ ) were both influenced. For the  $d'$ , we found that the shapes  
598 associated with good person (2.46, 95% CI[2.21 2.72]) is greater than shapes tagged with  
599 moral bad (2.07, 95% CI[1.83 2.32]),  $P_{PosteriorComparison} = 1$ . Shape tagged with morally  
600 good person is also greater than shapes tagged with neutral person (2.23, 95% CI[1.95  
601 2.49]),  $P_{PosteriorComparison} = 0.97$ . Also, the shapes tagged with neutral person is greater  
602 than shapes tagged with morally bad person,  $P_{PosteriorComparison} = 0.92$ .

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

608 *Reaction times.*

609 We fitted a Bayesian hierarchical GLM for RTs, with a log-normal distribution as the  
610 link function. We used the posterior distribution of the regression coefficient to make  
611 statistical inferences. As in previous studies, the matched conditions are much faster than  
612 the mismatched trials ( $P_{PosteriorComparison} = 1$ ). We focused on matched trials only, and

613 compared different conditions: Good is faster than the neutral,  $P_{PosteriorComparison} = .99$ ,  
614 it was also faster than the Bad condition,  $P_{PosteriorComparison} = 1$ . And the neutral  
615 condition is faster than the bad condition,  $P_{PosteriorComparison} = .99$ . However, the  
616 mismatched trials are largely overlapped. See Figure ??.

617 **HDDM.**

618 We fitted our data with HDDM, using the response-coding (See also, Hu et al., 2020).  
619 We estimated separate drift rate ( $v$ ), non-decision time ( $T_0$ ), and boundary separation ( $a$ )  
620 for each condition. We found that the shapes tagged with good person has higher drift rate  
621 and higher boundary separation than shapes tagged with both neutral and bad person.  
622 Also, the shapes tagged with neutral person has a higher drift rate than shapes tagged  
623 with bad person, but not for the boundary separation. Finally, we found that shapes  
624 tagged with bad person had longer non-decision time (see Figure ??).

625 **Experiment 1b**

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

628 **Method.**

629 **Participants.**

630 72 college students (49 female, age =  $20.17 \pm 2.08$  years) participated. 39 of them  
631 were recruited from Tsinghua University community in 2014; 33 were recruited from  
632 Wenzhou University in 2017. All participants were right-handed except one, and all had  
633 normal or corrected-to-normal vision. Informed consent was obtained from all participants  
634 prior to the experiment according to procedures approved by the local ethics committees.  
635 20 participant's data were excluded from analysis because nearly random level of accuracy,  
636 leaving 52 participants (36 female, age =  $20.25 \pm 2.31$  years).

637 **Stimuli and Tasks.** Three geometric shapes (triangle, square, and circle, with  $3.7^\circ$

638  $\times 3.7^\circ$  of visual angle) were presented above a white fixation cross subtending  $0.8^\circ \times 0.8^\circ$

639 of visual angle at the center of the screen. The three shapes were randomly assigned to

640 three labels with different moral valence: a morally bad person (" ", ERen), a morally

641 good person (" ", ShanRen) or a morally neutral person (" ", ChangRen). The order of

642 the associations between shapes and labels was counterbalanced across participants. Three

643 labels used in this experiment is selected based on the rating results from an independent

644 survey, in which participants rated the familiarity, frequency, and concreteness of eight

645 different words online. Of the eight words, three of them are morally positive (HaoRen,

646 ShanRen, Junzi), two of them are morally neutral (ChangRen, FanRen), and three of them

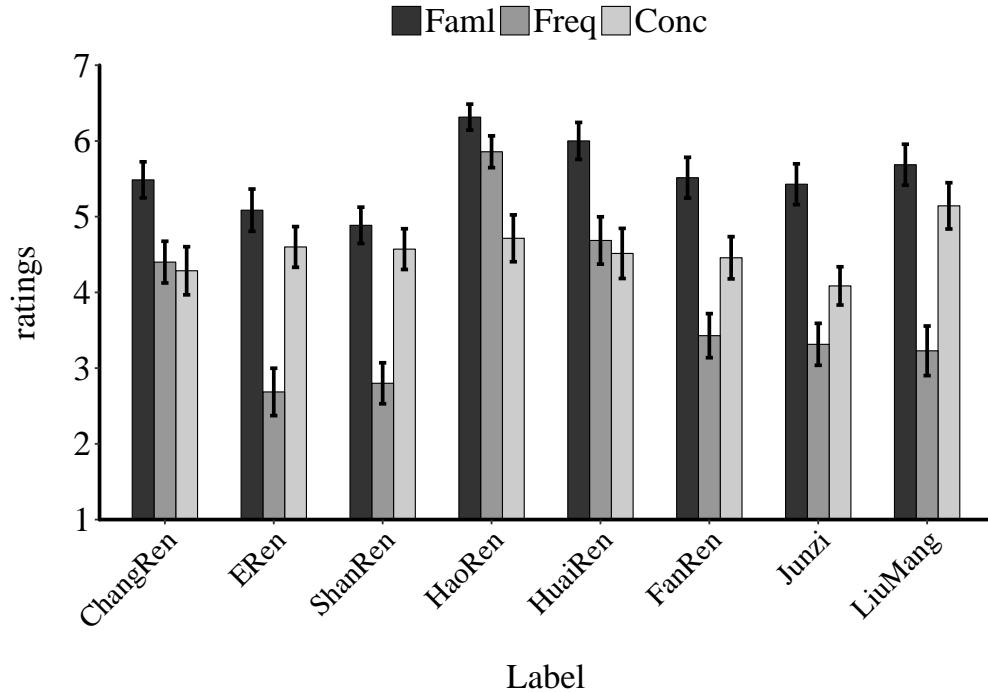
647 are morally negative (HuaiRen, ERen, LiuMang). An independent sample consist of 35

648 participants (22 females, age  $20.6 \pm 3.11$ ) were recruited to rate these words. Based on the

649 ratings (see supplementary materials Figure S1), we selected ShanRen, ChangRen, and

650 ERen to represent morally positive, neutral, and negative person.

#### Ratings for each label



651 **Procedure.**

652

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

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

656 **Results.**

657 **NHST.**

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

659  $d$  prime.

660 Repeated measures ANOVA revealed main effect of valence,  $F(1.83, 93.20) = 14.98$ ,  
661  $MSE = 0.18$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .053$ . Paired t test showed that the Good-Person condition  
662 ( $1.87 \pm 0.102$ ) was with greater  $d$  prime than Neutral condition ( $1.44 \pm 0.101$ ,  $t(51) =$   
663  $5.945$ ,  $p < 0.001$ ). We also found that the Bad-Person condition ( $1.67 \pm 0.11$ ) has also  
664 greater  $d$  prime than neutral condition ,  $t(51) = 3.132$ ,  $p = 0.008$ ). There Good-person  
665 condition was also slightly greater than the bad condition,  $t(51) = 2.265$ ,  $p = 0.0701$ .

666 *Reaction times.*

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

676 **BGLM.**

677        *Signal detection theory analysis of accuracy.*

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

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

691        *Reaction time.*

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

700        **HDDM.**

701        We found that the shapes tagged with good person has higher drift rate and higher  
702        boundary separation than shapes tagged with both neutral and bad person. Also, the

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

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

711 **Experiment 1c**

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

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

716 **Method.**

717 ***Participants.***

718 23 college students (15 female, age =  $22.61 \pm 2.62$  years) participated. All of them  
719 were recruited from Tsinghua University community in 2014. Informed consent was  
720 obtained from all participants prior to the experiment according to procedures approved by  
721 the local ethics committees. No participant was excluded because they overall accuracy  
722 were above 0.6.

723 ***Stimuli and Tasks.***

724 Three geometric shapes (triangle, square, and circle, with  $3.7^\circ \times 3.7^\circ$  of visual angle)  
725 were presented above a white fixation cross subtending  $0.8^\circ \times 0.8^\circ$  of visual angle at the  
726 center of the screen. The three most common names were chosen, which are neutral in  
727 moral valence before the manipulation. Three names (Zhang, Wang, Li) were first paired

728 with three paragraphs of behavioral description. Each description includes one sentence of  
729 biographic information and four sentences that describing the moral behavioral under that  
730 name. To assess the that these three descriptions represented good, neutral, and bad  
731 valence, we collected the ratings of three person on six dimensions: morality, likability,  
732 trustworthiness, dominance, competence, and aggressiveness, from an independent sample  
733 ( $n = 34$ , 18 female, age =  $19.6 \pm 2.05$ ). The rating results showed that the person with  
734 morally good behavioral description has higher score on morality ( $M = 3.59$ ,  $SD = 0.66$ )  
735 than neutral ( $M = 0.88$ ,  $SD = 1.1$ ),  $t(33) = 12.94$ ,  $p < .001$ , and bad conditions ( $M = -3.4$ ,  
736  $SD = 1.1$ ),  $t(33) = 30.78$ ,  $p < .001$ . Neutral condition was also significant higher than bad  
737 conditions  $t(33) = 13.9$ ,  $p < .001$  (See supplementary materials).

738 ***Procedure.***

739 After arriving the lab, participants were informed to complete two experimental  
740 tasks, first a social memory task to remember three person and their behaviors, after tested  
741 for their memory, they will finish a perceptual matching task. In the social memory task,  
742 the descriptions of three person were presented without time limitation. Participant  
743 self-paced to memorized the behaviors of each person. After they memorizing, a  
744 recognition task was used to test their memory effect. Each participant was required to  
745 have over 95% accuracy before preceding to matching task. The perceptual learning task  
746 was followed, three names were randomly paired with geometric shapes. Participants were  
747 required to learn the association and perform a practicing task before they start the formal  
748 experimental blocks. They kept practicing until they reached 70% accuracy. Then, they  
749 would start the perceptual matching task as in experiment 1a. They finished 6 blocks of  
750 perceptual matching trials, each have 120 trials.

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

752 **Results.** Figure ?? shows  $d'$  prime and reaction times of experiment 1c. We  
753 conducted same analysis as in Experiment 1a. Our analysis didn't show effect of valence

<sup>754</sup> on  $d$  prime,  $F(1.93, 42.56) = 0.23$ ,  $MSE = 0.41$ ,  $p = .791$ ,  $\hat{\eta}_G^2 = .005$ . Neither the effect of  
<sup>755</sup> valence on RT ( $F(1.63, 35.81) = 0.22$ ,  $MSE = 2,212.71$ ,  $p = .761$ ,  $\hat{\eta}_G^2 = .001$ ) or  
<sup>756</sup> interaction between valence and matchness on RT ( $F(1.79, 39.43) = 1.20$ ,  
<sup>757</sup>  $MSE = 1,973.91$ ,  $p = .308$ ,  $\hat{\eta}_G^2 = .005$ ).

<sup>758</sup> ***Signal detection theory analysis of accuracy.***

<sup>759</sup> We fitted a Bayesian hierarchical GLM for SDT. The results showed that when the  
<sup>760</sup> shapes were tagged with labels with different moral valence, the sensitivity ( $d'$ ) and criteria  
<sup>761</sup> ( $c$ ) were both influenced. For the  $d'$ , we found that the shapes tagged with morally good  
<sup>762</sup> person (2.30, 95% CI[1.93 2.70]) is greater than shapes tagged with moral bad (2.11, 95%  
<sup>763</sup> CI[1.83 2.42]),  $P_{PosteriorComparison} = 0.8$ . Shape tagged with morally good person is also  
<sup>764</sup> greater than shapes tagged with neutral person (2.16, 95% CI[1.88 2.45]),  
<sup>765</sup>  $P_{PosteriorComparison} = 0.75$ .

<sup>766</sup> Interesting, we also found the criteria for three conditions also differ, the shapes  
<sup>767</sup> tagged with good person has the highest criteria (-0.97, [-1.12 -0.82]), followed by shapes  
<sup>768</sup> tagged with neutral person(-0.96, [-1.09 -0.83]), and then the shapes tagged with bad  
<sup>769</sup> person(-1.03, [-1.22 -0.84]). However, pair-wise comparison showed that only showed strong  
<sup>770</sup> evidence for the difference between good and bad conditions.

<sup>771</sup> ***Reaction time.***

<sup>772</sup> We fitted a Bayesian hierarchical GLM for RTs, with a log-normal distribution as the  
<sup>773</sup> link function. We used the posterior distribution of the regression coefficient to make  
<sup>774</sup> statistical inferences. As in previous studies, the matched conditions are much faster than  
<sup>775</sup> the mismatched trials ( $P_{PosteriorComparison} = .75$ ). We focused on matched trials only, and  
<sup>776</sup> compared different conditions: Good () is not faster than the neutral (),  
<sup>777</sup>  $P_{PosteriorComparison} = .5$ , it was faster than the Bad condition (),  
<sup>778</sup>  $P_{PosteriorComparison} = .88$ . And the neutral condition is faster than the bad condition,  
<sup>779</sup>  $P_{PosteriorComparison} = .95$ . However, the mismatched trials are largely overlapped.

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

781 al., 2020). We estimated separate drift rate ( $v$ ), non-decision time ( $T_0$ ), and boundary

782 separation ( $a$ ) for each condition. We found that the shapes tagged with good person has

783 higher drift rate and higher boundary separation than shapes tagged with both neutral and

784 bad person. Also, the shapes tagged with neutral person has a higher drift rate than

785 shapes tagged with bad person, but not for the boundary separation. Finally, we found

786 that shapes tagged with bad person had longer non-decision time (see figure ??)).

787       **Experiment 2: Sequential presenting**

788       Experiment 2 was conducted for two purpose: (1) to further confirm the facilitation

789 effect of positive moral associations; (2) to test the effect of expectation of occurrence of

790 each pair. In this experiment, after participant learned the association between labels and

791 shapes, they were presented a label first and then a shape, they then asked to judge

792 whether the shape matched the label or not (see (Sui, Sun, Peng, & Humphreys, 2014)).

793 Previous studies showed that when the labels presented before the shapes, participants

794 formed expectations about the shape, and therefore a top-down process were introduced

795 into the perceptual matching processing. If the facilitation effect of positive moral valence

796 we found in experiment 1 was mainly drive by top-down processes, this sequential

797 presenting paradigm may eliminate or attenuate this effect; if, however, the facilitation

798 effect occurred because of button-up processes, then, similar facilitation effect will appear

799 even with sequential presenting paradigm.

800       **Method.**

801       **Participants.**

802       35 participants (17 female, age =  $21.66 \pm 3.03$ ) were recruited. 24 of them had

803 participated in Experiment 1a (9 male, mean age = 21.9, s.d. = 2.9), and the time gap

804 between these experiment 1a and experiment 2 is at least six weeks. The results of 1

805 participants were excluded from analysis because of less than 60% overall accuracy,  
806 remains 34 participants (17 female, age =  $21.74 \pm 3.04$ ).

807 ***Procedure.***

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

821 ***Data analysis.***

822 Data was analyzed as in study 1a.

823 **Results.**

824 ***NHST.***

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

827  ***$d'$  prime.***

828 There was evidence for the main effect of valence,  $F(1.83, 60.36) = 14.41$ ,  
829  $MSE = 0.23$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .066$ . Paired t test showed that the Good-Person condition

830 (2.79 ± 0.17) was with greater  $d$  prime than Netural condition (2.21 ± 0.16,  $t(33) = 4.723$ ,  
 831  $p = 0.001$ ) and Bad-person condition (2.41 ± 0.14),  $t(33) = 4.067$ ,  $p = 0.008$ ). There was  
 832 no-significant difference between Neutral-person and Bad-person conidition,  $t(33) = -1.802$ ,  
 833  $p = 0.185$ .

834 *Reaction time.*

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

837 We found interaction between Matchness and Valence ( $F(1.99, 65.70) = 9.53$ ,  
 838  $MSE = 605.36$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .017$ ) and then analyzed the matched trials and  
 839 mismatched trials separately, as in experiment 1a. For matched trials, we found the effect  
 840 of valence  $F(1.99, 65.76) = 10.57$ ,  $MSE = 1,192.65$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .067$ . Post-hoc  $t$ -tests  
 841 revealed that shapes associated with Good Person (548 ± 9.4) were responded faster than  
 842 Neutral-Person (582 ± 10.9), ( $t(33) = -3.95$ ,  $p = 0.0011$ ) and Bad Person (582 ± 10.2),  
 843  $t(33) = -3.9$ ,  $p = 0.0013$ ). While there was no significant differences between Neutral and  
 844 Bad-Person condition ( $t(33) = -0.01$ ,  $p = 0.999$ ). For non-matched trials, there was no  
 845 significant effect of Valence ( $F(1.99, 65.83) = 0.17$ ,  $MSE = 489.80$ ,  $p = .843$ ,  $\hat{\eta}_G^2 = .001$ ).

846 **BGLMM.**

847 *Signal detection theory analysis of accuracy.*

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

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

861 *Reaction times.*

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

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

878 **Discussion**

879 In this experiment, we repeated the results pattern that the positive moral valenced  
880 stimuli has an advantage over the neutral or the negative valence association. Moreover,

881 with a cross-task analysis, we did not find evidence that the experiment task interacted  
882 with moral valence, suggesting that the effect might not be effect by experiment task.  
883 These findings suggested that the facilitation effect of positive moral valence is robust and  
884 not affected by task. This robust effect detected by the associative learning is unexpected.

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

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

888 **Method.**

889 **Participants.**

890 24 college students (8 female, age =  $22.88 \pm 2.79$ ) participated the current study, all  
891 of them were from Tsinghua University in 2014. Informed consent was obtained from all  
892 participants prior to the experiment according to procedures approved by a local ethics  
893 committee. No participant was excluded from behavioral analysis.

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

898 **Stimuli.**

899 Three geometric shapes (triangle, square and circle, each  $4.6^\circ \times 4.6^\circ$  of visual angle)  
900 were presented at the center of screen for 50 ms after 500ms of fixation ( $0.8^\circ \times 0.8^\circ$  of  
901 visual angle). The association of the three shapes to bad person (" , HuaiRen"), good  
902 person (" , HaoRen") or ordinary person (" , ChangRen") was counterbalanced across  
903 participants. The words bad person, good person or ordinary person ( $3.6^\circ \times 1.6^\circ$ ) was also  
904 displayed at the center fo the screen. Participants had to judge whether the pairings of  
905 label and shape matched (e.g., Does the circle represent a bad person?). The experiment

906 was run on a PC using E-prime software (version 2.0). These stimuli were displayed on a  
907 22-in CRT monitor ( $1024 \times 768$  at 100Hz). We used backward masking to avoid  
908 over-processing of the moral words, in which a scrambled picture were presented for 900 ms  
909 after the label. Also, to avoid the ceiling effect on accuracy, shapes were presented on a  
910 noisy background based on our pilot studies. The noisy images were made by scrambling a  
911 picture of 3/4 gray and 1/4 white at resolution of  $2 \times 2$  pixel.

912 ***Procedure.***

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

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

924 All the stimuli were presented on a gray background (RGB: 127, 127, 127). E-primed  
925 2.0 was used to present stimuli and collect behavioral results. Data were collected and  
926 analyzed when accuracy performance in total reached 60%.

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

928 **Results.**

929 **NHST.**

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

932 *d prime.*

933 We conducted repeated measures ANOVA, with moral valence as independent  
 934 variable. The results revealed the main effect of valence ( $F(1.74, 40.05) = 3.76$ ,  
 935  $MSE = 0.10$ ,  $p = .037$ ,  $\hat{\eta}_G^2 = .021$ ). Post-hoc analysis revealed that shapes link with Good  
 936 person (mean = 3.13, SE = 0.109) is greater than Neutral condition (mean = 2.88, SE =  
 937 0.14),  $t = 2.916$ ,  $df = 24$ ,  $p = 0.02$ , p-value adjusted by Tukey method, but the *d* prime  
 938 between Good and bad (mean = 3.03, SE = 0.142) ( $t = 1.512$ ,  $df = 24$ ,  $p = 0.3034$ , p-value  
 939 adjusted by Tukey method), bad and neutral ( $t = 1.599$ ,  $df = 24$ ,  $p = 0.2655$ , p-value  
 940 adjusted by Tukey method) were not significant.

941 *Reaction times.*

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

944 We found intercation between Matchness and Valence ( $F(1.97, 45.20) = 20.45$ ,  
 945  $MSE = 450.47$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .021$ ) and then analyzed the matched trials and  
 946 mismatched trials separately, as in experiment 2. For matched trials, we found the effect of  
 947 valence  $F(1.97, 45.25) = 32.37$ ,  $MSE = 522.42$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .078$ . For non-matched  
 948 trials, there was no significant effect of Valence ( $F(1.77, 40.67) = 0.35$ ,  $MSE = 242.15$ ,  
 949  $p = .679$ ,  $\hat{\eta}_G^2 = .000$ ). Post-hoc *t*-tests revealed that shapes associated with Good Person  
 950 (mean = 550, SE = 13.8) were responded faster than Neutral-Person (501, SE = 14.7),  
 951 ( $t(24) = -5.171$ ,  $p = 0.0001$ ) and Bad Person (523, SE = 16.3),  $t(24) = -8.137$ ,  $p <$   
 952 0.0001., and Neutral is faster than Bad-Person condition ( $t(32) = -3.282$ ,  $p = 0.0085$ ).

953 **BGLM.**

954 *Signal detection theory analysis of accuracy.*

955 *Reaction time.*

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

958 separation (*a*) for each condition. We found that, similar to experiment 2, the shapes  
959 tagged with good person has higher drift rate and higher boundary separation than shapes  
960 tagged with both neutral and bad person, but only for the self-referential condition. Also,  
961 the shapes tagged with neutral person has a higher drift rate than shapes tagged with bad  
962 person, but not for the boundary separation, and this effect also exist only for the  
963 self-referential condition.

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

## 968 Part 2: interaction between valence and identity

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

### 972 Experiment 3a

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

#### 977 Method.

##### 978 Participants.

979 38 college students (15 female, age =  $21.92 \pm 2.16$ ) participated in experiment 3a.  
980 All of them were right-handed, and all had normal or corrected-to-normal vision. Informed  
981 consent was obtained from all participants prior to the experiment according to procedures

982 approved by a local ethics committee. One female and one male student did not finish the  
983 experiment, and 1 participants' data were excluded from analysis because less than 60%  
984 overall accuracy, remains 35 participants (13 female, age =  $22.11 \pm 2.13$ ).

985 ***Design.***

986 Study 3a combined moral valence with self-relevance, hence the experiment has a  $2 \times$   
987  $3 \times 2$  within-subject design. The first variable was self-relevance, include two levels:  
988 self-relevance vs. stranger-relevance; the second variable was moral valence, include good,  
989 neutral and bad; the third variable was the matching between shape and label: match  
990 vs. nonmatch.

991 ***Stimuli.***

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

997 ***Procedure.***

998 After being fully explained and signed the informed consent, participants were  
999 instructed to chose a name that can represent a stranger with same gender as the  
1000 participant themselves, from a common Chinese name pool. Before experiment, the  
1001 experimenter explained the meaning of each label to participants. For example, the "good  
1002 self" mean the morally good side of themselves, them could imagine the moment when they  
1003 do something's morally applauded, "bad self" means the morally bad side of themselves,  
1004 they could also imagine the moment when they doing something morally wrong, and  
1005 "neutral self" means the aspect of self that does not related to morality, they could imagine  
1006 the moment when they doing something irrelevant to morality. In the same sense, the  
1007 "good other", "bad other", and "neutral other" means the three different aspects of the

stranger, whose name was chosen before the experiment. Then, the experiment proceeded as study 1a. Each participant finished 6 blocks, each have 120 trials. The sequence of trials was pseudo-randomized so that there are 10 matched trials for each condition and 10 non-matched trials for each condition (good self, neutral self, bad self, good other, neutral other, bad other) for each block.

**1013      *Data Analysis.***

1014      Data analysis followed strategies described in the general method section. Reaction  
1015     times and  $d$  prime data were analyzed as in study 1 and study 2, except that one more  
1016     within-subject variable (i.e., self-relevance) was included in the analysis.

1017      **Results.**

1018      **NHST.**

1019      Figure 2 shows  $d$  prime and reaction times of experiment 3a. Less than 5% correct  
1020     trials with less than 200ms reaction times were excluded.

1021       $d$  prime.

1022      There was evidence for the main effect of valence,  $F(1.89, 64.37) = 11.09$ ,  
1023      $MSE = 0.23$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .039$ , and main effect of self-relevance,  $F(1, 34) = 3.22$ ,  
1024      $MSE = 0.54$ ,  $p = .082$ ,  $\hat{\eta}_G^2 = .015$ , as well as the interaction,  $F(1.79, 60.79) = 3.39$ ,  
1025      $MSE = 0.43$ ,  $p = .045$ ,  $\hat{\eta}_G^2 = .022$ .

1026      We then conducted separated ANOVA for self-referential and other-referential trials.  
1027      The valence effect was shown for the self-referential conditions,  $F(1.65, 56.25) = 13.98$ ,  
1028      $MSE = 0.31$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .119$ . Post-hoc test revealed that the Good-Self condition  
1029      $(1.97 \pm 0.14)$  was with greater  $d$  prime than Netural condition  $(1.41 \pm 0.12$ ,  $t(34) = 4.505$ ,  
1030      $p = 0.0002$ ), and Bad-self condition  $(1.43 \pm 0.102)$ ,  $t(34) = 3.856$ ,  $p = 0.0014$ . There was  
1031     difference between neutral and bad condition,  $t(34) = -0.238$ ,  $p = 0.9694$ . However, no  
1032     effect of valence was found for the other-referential condition  $F(1.98, 67.36) = 0.38$ ,  
1033      $MSE = 0.35$ ,  $p = .681$ ,  $\hat{\eta}_G^2 = .004$ .

1034       *Reaction time.*

1035       We found interaction between Matchness and Valence ( $F(1.98, 67.44) = 26.29$ ,

1036        $MSE = 730.09$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .025$ ) and then analyzed the matched trials and nonmatch

1037       trials separately, as in previous experiments.

1038       For the match trials, we found that the interaction between identity and valence,

1039        $F(1.72, 58.61) = 3.89$ ,  $MSE = 2,750.19$ ,  $p = .032$ ,  $\hat{\eta}_G^2 = .019$ , as well as the main effect of

1040       valence  $F(1.98, 67.34) = 35.76$ ,  $MSE = 1,127.25$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .079$ , but not the effect of

1041       identity  $F(1, 34) = 0.20$ ,  $MSE = 3,507.14$ ,  $p = .660$ ,  $\hat{\eta}_G^2 = .001$ . As for the  $d$  prime, we

1042       separated analyzed the self-referential and other-referential trials. For the Self-referential

1043       trials, we found the main effect of valence,  $F(1.80, 61.09) = 30.39$ ,  $MSE = 1,584.53$ ,

1044        $p < .001$ ,  $\hat{\eta}_G^2 = .159$ ; for the other-referential trials, the effect of valence is weaker,

1045        $F(1.86, 63.08) = 2.85$ ,  $MSE = 2,224.30$ ,  $p = .069$ ,  $\hat{\eta}_G^2 = .024$ . We then focused on the self

1046       conditions: the good-self condition ( $713 \pm 12$ ) is faster than neutral- ( $776 \pm 11.8$ ),  $t(34) =$

1047        $-7.396$ ,  $p < .0001$ , and bad-self ( $772 \pm 10.1$ ) conditions,  $t(34) = -5.66$ ,  $p < .0001$ . But

1048       there is not difference between neutral- and bad-self conditions,  $t(34) = 0.481$ ,  $p = 0.881$ .

1049       For the nonmatch trials, we didn't found any strong effect: identity,  $F(1, 34) = 3.43$ ,

1050        $MSE = 660.02$ ,  $p = .073$ ,  $\hat{\eta}_G^2 = .004$ , valence  $F(1.89, 64.33) = 0.40$ ,  $MSE = 444.10$ ,

1051        $p = .661$ ,  $\hat{\eta}_G^2 = .001$ , or interaction between the two  $F(1.94, 66.02) = 2.42$ ,  $MSE = 817.35$ ,

1052        $p = .099$ ,  $\hat{\eta}_G^2 = .007$ .

1053       **BGLM.**

1054       *Signal detection theory analysis of accuracy.*

1055       We found that the  $d$  prime is greater when shapes were associated with good self

1056       condition than with neutral self or bad self, but shapes associated with bad self and neutral

1057       self didn't show differences. Comparing the self vs other under three condition revealed

1058       that shapes associated with good self is greater than with good other, but with a weak

1059       evidence. In contrast, for both neutral and bad valence condition, shapes associated with

1060 other had greater  $d$  prime than with self.

1061 *Reaction time.*

1062 In reaction times, we found that same trends in the match trials as in the RT: while  
1063 the shapes associated with good self was greater than with good other (log mean diff =  
1064 -0.02858, 95%HPD[-0.070898, 0.0154]), the direction is reversed for neutral and negative  
1065 condition. see Figure 3

1066 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et  
1067 al., 2020). We estimated separate drift rate ( $v$ ), non-decision time ( $T_0$ ), and boundary  
1068 separation ( $a$ ) for each condition. We found that the shapes tagged with good person has  
1069 higher drift rate and higher boundary separation than shapes tagged with both neutral and  
1070 bad person. Also, the shapes tagged with neutral person has a higher drift rate than  
1071 shapes tagged with bad person, but not for the boundary separation. Finally, we found  
1072 that shapes tagged with bad person had longer non-decision time (see figure 4)).

1073 **Experiment 3b**

1074 In study 3a, participants had to remember 6 pairs of association, which cause high  
1075 cognitive load during the whole experiment. To eliminate the influence of cognitive load, we  
1076 conducted study 3b, in which participant learn three aspect of self and stranger separately  
1077 in to consecutive task. We hypothesize that we will replicate the pattern of study 3a, i.e.,  
1078 the effect of moral valence only occurs for self-relevant conditions. ### Method

1079 **Participants.**

1080 Study 3b were finished in 2017, at that time we have calculated that the effect size  
1081 (Cohen's  $d$ ) of good-person (or good-self) vs. bad-person (or bad-other) was between 0.47 ~  
1082 0.53, based on the data from Tsinghua community in study 1a, 1b, 2, 3a, 4a, and 4b. Based  
1083 on this effect size, we estimated that 54 participants would allow we to detect the effect  
1084 size of Cohen's  $= 0.5$  with 95% power and alpha = 0.05, using G\*power 3.192 (Faul,

1085 Erdfelder, Buchner, & Lang, 2009). Therefore, we planned to stop after we arrived this  
1086 number. During the data collected at Wenzhou University, 61 participants (45 females; 19  
1087 to 25 years of age, age =  $20.42 \pm 1.77$ ) came to the testing room and we tested all of them  
1088 during a single day. All participants were right-handed, and all had normal or  
1089 corrected-to-normal vision. Informed consent was obtained from all participants prior to  
1090 the experiment according to procedures approved by a local ethics committee. 4  
1091 participants' data were excluded from analysis because their over all accuracy was lower  
1092 than 60%, 1 more participant was excluded because of zero hit rate for one condition,  
1093 leaving 56 participants (43 females; 19 to 25 years old, age =  $20.27 \pm 1.60$ ).

1094        ***Design.***

1095        Study 3b has the same experimental design as 3a, with a  $2 \times 3 \times 2$  within-subject  
1096 design. The first variable was self-relevance, include two levels: self-relevant  
1097 vs. stranger-relevant; the second variable was moral valence, include good, neutral and bad;  
1098 the third variable was the matching between shape and label: match vs. mismatch.  
1099 Stimuli. The stimuli used in study 3b share the same parameters with experiment 3a. 6  
1100 shapes were included (triangle, square, circle, trapezoid, diamond, regular pentagon), as  
1101 well as 6 labels, but the labels changed to “good self”, “neutral self”, “bad self”, “good  
1102 him/her”, “bad him/her”, “neutral him/her”, the stranger's label is consistent with  
1103 participants' gender. Same as study 3a, we asked participant to chosen an unfamiliar name  
1104 of their own gender to be the stranger before showing them the relationship. Note, because  
1105 of implementing error, the personal distance data did not collect for this experiment.

1106        ***Stimuli.***

1107        The stimuli used in study 3b is the same as in experiment 3a.

1108        ***Procedure.***

1109        In this experiment, participants finished two matching tasks, i.e., self-matching task,  
1110 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.

**1120      *Data Analysis.***

1121      Same as experiment 3a.

1122      **Results.**

1123      **NHST.**

1124      Figure 5 shows  $d$  prime and reaction times of experiment 3b. Less than 5% correct trials with less than 200ms reaction times were excluded.

1126       $d$  prime.

1127      There was no evidence for the main effect of valence,  $F(1.92, 105.43) = 1.90$ ,  
1128       $MSE = 0.33$ ,  $p = .157$ ,  $\hat{\eta}_G^2 = .005$ , but we found a main effect of self-relevance,  
1129       $F(1, 55) = 4.65$ ,  $MSE = 0.89$ ,  $p = .035$ ,  $\hat{\eta}_G^2 = .017$ , as well as the interaction,  
1130       $F(1.90, 104.36) = 5.58$ ,  $MSE = 0.26$ ,  $p = .006$ ,  $\hat{\eta}_G^2 = .011$ .

1131      We then conducted separated ANOVA for self-referential and other-referential trials.  
1132      The valence effect was shown for the self-referential conditions,  $F(1.75, 96.42) = 6.73$ ,  
1133       $MSE = 0.30$ ,  $p = .003$ ,  $\hat{\eta}_G^2 = .037$ . Post-hoc test revealed that the Good-Self condition  
1134      ( $2.15 \pm 0.12$ ) was with greater  $d$  prime than Neutral condition ( $1.83 \pm 0.12$ ,  $t(34) = 3.36$ ,  
1135       $p = 0.0031$ ), and Bad-self condition ( $1.87 \pm 0.12$ ),  $t(34) = 2.955$ ,  $p = 0.01$ . There was

<sub>1136</sub> difference between neutral and bad condition,  $t(34) = -0.039$ ,  $p = 0.914$ . However, no effect  
<sub>1137</sub> of valence was found for the other-referential condition  $F(1.93, 105.97) = 0.61$ ,  
<sub>1138</sub>  $MSE = 0.31$ ,  $p = .539$ ,  $\hat{\eta}_G^2 = .002$ .

<sub>1139</sub> *Reaction time.*

<sub>1140</sub> We found interaction between Matchness and Valence ( $F(1.86, 102.47) = 15.44$ ,  
<sub>1141</sub>  $MSE = 3, 112.78$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .006$ ) and then analyzed the matched trials and  
<sub>1142</sub> nonmatch trials separately, as in previous experiments.

<sub>1143</sub> For the match trials, we found that the interaction between identity and valence,  
<sub>1144</sub>  $F(1.67, 92.11) = 6.14$ ,  $MSE = 6, 472.48$ ,  $p = .005$ ,  $\hat{\eta}_G^2 = .009$ , as well as the main effect of  
<sub>1145</sub> valence  $F(1.88, 103.65) = 24.25$ ,  $MSE = 5, 994.25$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .038$ , but not the effect  
<sub>1146</sub> of identity  $F(1, 55) = 48.49$ ,  $MSE = 25, 892.59$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .153$ . As for the  $d$  prime,  
<sub>1147</sub> we separated analyzed the self-referential and other-referential trials. For the  
<sub>1148</sub> Self-referential trials, we found the main effect of valence,  $F(1.66, 91.38) = 23.98$ ,  
<sub>1149</sub>  $MSE = 6, 965.61$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .100$ ; for the other-referential trials, the effect of valence  
<sub>1150</sub> is weaker,  $F(1.89, 103.94) = 5.96$ ,  $MSE = 5, 589.90$ ,  $p = .004$ ,  $\hat{\eta}_G^2 = .014$ . We then focused  
<sub>1151</sub> on the self conditions: the good-self condition ( $713 \pm 12$ ) is faster than neutral- ( $776 \pm$   
<sub>1152</sub>  $11.8$ ),  $t(34) = -7.396$ ,  $p < .0001$ , and bad-self ( $772 \pm 10.1$ ) conditions,  $t(34) = -5.66$ ,  $p <$   
<sub>1153</sub>  $.0001$ . But there is not difference between neutral- and bad-self conditions,  $t(34) = 0.481$ ,  $p$   
<sub>1154</sub>  $= 0.881$ .

<sub>1155</sub> For the nonmatch trials, we didn't found any strong effect: identity,  $F(1, 55) = 10.31$ ,  
<sub>1156</sub>  $MSE = 24, 590.52$ ,  $p = .002$ ,  $\hat{\eta}_G^2 = .035$ , valence  $F(1.98, 108.63) = 20.57$ ,  $MSE = 2, 847.51$ ,  
<sub>1157</sub>  $p < .001$ ,  $\hat{\eta}_G^2 = .016$ , or interaction between the two  $F(1.93, 106.25) = 35.51$ ,  
<sub>1158</sub>  $MSE = 1, 939.88$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .019$ .

<sub>1159</sub> **BGLM.**

<sub>1160</sub> *Signal detection theory analysis of accuracy.*

<sub>1161</sub> 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 6

**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 7)).

## Experiment 6b

Experiment 6b was conducted to study the neural correlates of the prioritization effect of positive self, i.e., the neural underlying of the behavioral effect found int

1187 experiment 3a. However, as in experiment 6a, the procedure of this experiment was  
1188 modified to adopted to ERP experiment.

1189       **Method.**

1190       ***Participants.***

1191       23 college students (8 female, age =  $22.86 \pm 2.47$ ) participated the current study, all  
1192 of them were recruited from Tsinghua University in 2016. Informed consent was obtained  
1193 from all participants prior to the experiment according to procedures approved by a local  
1194 ethics committee. For day 1's data, 1 participant was excluded from the current analysis  
1195 because of lower than 60% overall accuracy, remaining 22 participants (8 female, age =  
1196  $22.76 \pm 2.49$ ). For day 2's data, one participant dropped out, leaving 22 participants (9  
1197 female, age =  $23.05 \pm 2.46$ ), all of them has overall accuracy higher than 60%.

1198       ***Design.***

1199       The experimental design of this experiment is same as experiment 3: a  $2 \times 3 \times 2$   
1200 within-subject design with self-relevance (self-relevant vs. other-relevant), moral valence  
1201 (good, neutral, and bad) and matchness between shape and label (match vs. mismatch) as  
1202 within-subject variables.

1203       ***Stimuli.***

1204       As in experiment 3a, 6 shapes were included (triangle, square, circle, trapezoid,  
1205 diamond, regular pentagon), as well as 6 labels (good self, neutral self, bad self, good  
1206 person, bad person, neutral person). To match the concreteness of the label, we asked  
1207 participant to chosen an unfamiliar name of their own gender to be the stranger.

1208       ***Procedure.***

1209       The procedure was similar to Experiment 2 and 6a. Subjects first learned the  
1210 associations between labels and shapes and then completed a shape-label matching task. In  
1211 each trial of the matching task, a fixation were first presented for 500 ms, followed by a 50

1212 ms label; then, a scrambled picture presented 900 ms. After the backward mask, the shape  
1213 were presented on a noisy background for 50ms. Participant have to response in 1000ms  
1214 after the presentation of the shape, and finally, a feedback screen was presented for 500 ms.  
1215 The inter-trial interval (ITI) were randomly varied at the range of 1000 ~ 1400 ms.

1216 All the stimuli were presented on a gray background (RGB: 127, 127, 127). E-primed  
1217 2.0 was used to present stimuli and collect behavioral results. Data were collected and  
1218 analyzed when accuracy performance in total reached 60%.

1219 Because learning 6 associations was more difficult than 3 associations and participant  
1220 might have low accuracy (see experiment 3a), the current study had extended to a two-day  
1221 paradigm to maximizing the accurate trials that can be used in EEG data. At the first day,  
1222 participants learnt the associations and finished 9 blocks of the matching task, each had  
1223 120 trials, without EEG recording. That is, each condition has 90 trials.

1224 Participants came back to lab at the second day and finish the same task again, with  
1225 EEG recorded. Before the EEG experiment, each participant finished a practice session  
1226 again, if their accuracy is equal or higher than 85%, they start the experiment (one  
1227 participant used lower threshold 75%). Each participant finished 18 blocks, each has 90  
1228 trials. One participant finished additional 6 blocks because of high error rate at the  
1229 beginning, another two participant finished addition 3 blocks because of the technique  
1230 failure in recording the EEG data. To increase the number of trials that can be used for  
1231 EEG data analysis, matched trials has twice number as mismatched trials, therefore, for  
1232 matched trials each participants finished 180 trials for each condition, for mismatched  
1233 trials, each conditions has 90 trials.

1234 ***Data Analysis.***

1235 Same as experiment 3a.

1236 **Results of Day 1.**

1237 **NHST.**

1238 Figure 8 shows *d* prime and reaction times of experiment 3b. Less than 5% correct

1239 trials with less than 200ms reaction times were excluded.

1240 *d prime.*

1241 There was no evidence for the main effect of valence,  $F(1.91, 40.20) = 11.98$ ,

1242  $MSE = 0.15$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .040$ , but we found a main effect of self-relevance,

1243  $F(1, 21) = 1.21$ ,  $MSE = 0.20$ ,  $p = .284$ ,  $\hat{\eta}_G^2 = .003$ , as well as the interaction,

1244  $F(1.28, 26.90) = 12.88$ ,  $MSE = 0.21$ ,  $p = .001$ ,  $\hat{\eta}_G^2 = .041$ .

1245 We then conducted separated ANOVA for self-referential and other-referential trials.

1246 The valence effect was shown for the self-referential conditions,  $F(1.73, 36.42) = 29.31$ ,

1247  $MSE = 0.14$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .147$ . Post-hoc test revealed that the Good-Self condition

1248 ( $2.15 \pm 0.12$ ) was with greater *d* prime than Neutral condition ( $1.83 \pm 0.12$ ,  $t(34) = 3.36$ ,

1249  $p = 0.0031$ ), and Bad-self condition ( $1.87 \pm 0.12$ ),  $t(34) = 2.955$ ,  $p = 0.01$ . There was

1250 difference between neutral and bad condition,  $t(34) = -0.039$ ,  $p = 0.914$ . However, no effect

1251 of valence was found for the other-referential condition  $F(1.75, 36.72) = 0.00$ ,  $MSE = 0.18$ ,

1252  $p = .999$ ,  $\hat{\eta}_G^2 = .000$ .

1253 *Reaction time.*

1254 We found interaction between Matchness and Valence ( $F(1.79, 37.63) = 4.07$ ,

1255  $MSE = 704.90$ ,  $p = .029$ ,  $\hat{\eta}_G^2 = .003$ ) and then analyzed the matched trials and nonmatch

1256 trials separately, as in previous experiments.

1257 For the match trials, we found that the interaction between identity and valence,

1258  $F(1.72, 36.16) = 4.55$ ,  $MSE = 1,560.90$ ,  $p = .022$ ,  $\hat{\eta}_G^2 = .015$ , as well as the main effect of

1259 valence  $F(1.93, 40.55) = 9.83$ ,  $MSE = 1,951.84$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .044$ , but not the effect of

1260 identity  $F(1, 21) = 4.87$ ,  $MSE = 2,032.05$ ,  $p = .039$ ,  $\hat{\eta}_G^2 = .012$ . As for the *d* prime, we

1261 separated analyzed the self-referential and other-referential trials. For the Self-referential

1262 trials, we found the main effect of valence,  $F(1.92, 40.38) = 14.48$ ,  $MSE = 1,647.20$ ,

1263  $p < .001$ ,  $\hat{\eta}_G^2 = .112$ ; for the other-referential trials, the effect of valence is weaker,

<sub>1264</sub>  $F(1.79, 37.50) = 1.04$ ,  $MSE = 1,842.07$ ,  $p = .356$ ,  $\hat{\eta}_G^2 = .008$ . We then focused on the self  
<sub>1265</sub> conditions: the good-self condition ( $713 \pm 12$ ) is faster than neutral- ( $776 \pm 11.8$ ),  $t(34) =$   
<sub>1266</sub>  $-7.396$ ,  $p < .0001$ , and bad-self ( $772 \pm 10.1$ ) conditions,  $t(34) = -5.66$ ,  $p < .0001$ . But  
<sub>1267</sub> there is not difference between neutral- and bad-self conditions,  $t(34) = 0.481$ ,  $p = 0.881$ .

<sub>1268</sub> For the nonmatch trials, we didn't found any strong effect: identity,  $F(1, 21) = 2.76$ ,  
<sub>1269</sub>  $MSE = 1,718.93$ ,  $p = .112$ ,  $\hat{\eta}_G^2 = .006$ , valence  $F(1.61, 33.77) = 3.81$ ,  $MSE = 1,532.21$ ,  
<sub>1270</sub>  $p = .041$ ,  $\hat{\eta}_G^2 = .012$ , or interaction between the two  $F(1.90, 39.97) = 2.23$ ,  $MSE = 720.80$ ,  
<sub>1271</sub>  $p = .123$ ,  $\hat{\eta}_G^2 = .004$ .

<sub>1272</sub> **BGLM.**

<sub>1273</sub> *Signal detection theory analysis of accuracy.*

<sub>1274</sub> We found that the  $d$  prime is greater when shapes were associated with good self  
<sub>1275</sub> condition than with neutral self or bad self, but shapes associated with bad self and neutral  
<sub>1276</sub> self didn't show differences. comparing the self vs other under three condition revealed that  
<sub>1277</sub> shapes associated with good self is greater than with good other, but with a weak evidence.  
<sub>1278</sub> In contrast, for both neutral and bad valence condition, shapes associated with other had  
<sub>1279</sub> greater  $d$  prime than with self.

<sub>1280</sub> *Reaction time.*

<sub>1281</sub> In reaction times, we found that same trends in the match trials as in the RT: while  
<sub>1282</sub> the shapes associated with good self was greater than with good other ( $\log$  mean diff =  
<sub>1283</sub>  $-0.02858$ , 95%HPD[-0.070898, 0.0154]), the direction is reversed for neutral and negative  
<sub>1284</sub> condition. see Figure 9

<sub>1285</sub> **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et  
<sub>1286</sub> al., 2020). We estimated separate drift rate ( $v$ ), non-decision time ( $T_0$ ), and boundary  
<sub>1287</sub> separation ( $a$ ) for each condition. We found that, similar to experiment 3a, the shapes  
<sub>1288</sub> tagged with good person has higher drift rate and higher boundary separation than shapes  
<sub>1289</sub> tagged with both neutral and bad person, but only for the self-referential condition. Also,

1290 the shapes tagged with neutral person has a higher drift rate than shapes tagged with bad  
1291 person, but not for the boundary separation, and this effect also exist only for the  
1292 self-referential condition.

1293 Interestingly, we found that in both self-referential and other-referential conditions,  
1294 the shapes associated bad valence have higher drift rate and higher boundary separation.  
1295 which might suggest that the shape associated with bad stimuli might be prioritized in the  
1296 non-match trials (see figure 10).

1297 **Part 3: Implicit binding between valence and identity**

1298 In this part, we reported two studies in which the moral valence or the self-referential  
1299 processing is not task-relevant. We are interested in testing whether the task-relevance will  
1300 eliminate the effect observed in previous experiment.

1301 **Experiment 4a: Morality as task-irrelevant variable**

1302 In part two (experiment 3a and 3b), participants learned the association between self  
1303 and moral valence directly. In Experiment 4a, we examined whether the interaction  
1304 between moral valence and identity occur even when one of the variable was irrelevant to  
1305 the task. In experiment 4a, participants learnt associations between shapes and self/other  
1306 labels, then made perceptual match judgments only about the self or other conditions  
1307 labels and shapes (cf. Sui et al. (2012)). However, we presented labels of different moral  
1308 valence in the shapes, which means that the moral valence factor become task irrelevant. If  
1309 the binding between moral good and self is intrinsic and automatic, then we will observe  
1310 that facilitating effect of moral good for self conditions, but not for other conditions.

1311 **Method.**

1312 ***Participants.***

1313        64 participants (37 female, age =  $19.70 \pm 1.22$ ) participated the current study, 32 of  
1314      them were from Tsinghua University in 2015, 32 were from Wenzhou University  
1315      participated in 2017. All participants were right-handed, and all had normal or  
1316      corrected-to-normal vision. Informed consent was obtained from all participants prior to  
1317      the experiment according to procedures approved by a local ethics committee. The data  
1318      from 5 participants from Wenzhou site were excluded from analysis because their accuracy  
1319      was close to chance ( $< 0.6$ ). The results for the remaining 59 participants (33 female, age  
1320      =  $19.78 \pm 1.20$ ) were analyzed and reported.

1321      ***Design.***

1322      As in Experiment 3, a  $2 \times 3 \times 2$  within-subject design was used. The first variable was  
1323      self-relevance (self and stranger associations); the second variable was moral valence (good,  
1324      neutral and bad associations); the third variable was the matching between shape and label  
1325      (matching vs. non-match for the personal association). However, in this the task,  
1326      participants only learn the association between two geometric shapes and two labels (self  
1327      and other), i.e., only self-relevance were related to the task. The moral valence  
1328      manipulation was achieved by embedding the personal label of the labels in the geometric  
1329      shapes, see below. For simplicity, the trials where shapes where paired with self and with a  
1330      word of “good person” inside were shorted as good-self condition, similarly, the trials where  
1331      shapes paired with the self and with a word of “bad person” inside were shorted as bad-self  
1332      condition. Hence, we also have six conditions: good-self, neutral-self, bad-self, good-other,  
1333      neutral-other, and bad-other.

1334      ***Stimuli.***

1335      2 shapes were included (circle, square) and each appeared above a central fixation  
1336      cross with the personal label appearing below. However, the shapes were not empty but  
1337      with a two-Chinese-character word in the middle, the word was one of three labels with  
1338      different moral valence: “good person”, “bad person” and “neutral person”. Before the

1339 experiment, participants learned the self/other association, and were informed to only  
1340 response to the association between shapes' configures and the labels below the fixation, but  
1341 ignore the words within shapes. Besides the behavioral experiments, participants from  
1342 Tsinghua community also finished questionnaires as Experiments 3, and participants from  
1343 Wenzhou community finished a series of questionnaire as the other experiment finished in  
1344 Wenzhou.

1345 ***Procedure.***

1346 The procedure was similar to Experiment 1. There were 6 blocks of trial, each with  
1347 120 trials for 2017 data. Due to procedure error, the data collected in 2015 in Tsinghua  
1348 community only have 60 trials for each block, i.e., 30 trials per condition.

1349 As in study 3a, before each task, the instruction showed the meaning of each label to  
1350 participants. The self-matching task and other-matching task were randomized between  
1351 participants. Each participant finished 6 blocks, each have 120 trials.

1352 ***Data Analysis.***

1353 Same as experiment 3a.

1354 **Results.**

1355 ***NHST.***

1356 Figure 11 shows  $d$  prime and reaction times of experiment 3a. Less than 5% correct  
1357 trials with less than 200ms reaction times were excluded.

1358  $d$  prime.

1359 There was no evidence for the main effect of valence,  $F(1.93, 111.66) = 0.53$ ,  
1360  $MSE = 0.12$ ,  $p = .581$ ,  $\hat{\eta}_G^2 = .000$ , but we found a main effect of self-relevance,  
1361  $F(1, 58) = 121.04$ ,  $MSE = 0.48$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .189$ , as well as the interaction,  
1362  $F(1.99, 115.20) = 4.12$ ,  $MSE = 0.14$ ,  $p = .019$ ,  $\hat{\eta}_G^2 = .004$ .

1363 We then conducted separated ANOVA for self-referential and other-referential trials.

1364 The valence effect was shown for the self-referential conditions,  $F(1.95, 112.92) = 3.01$ ,

1365  $MSE = 0.15$ ,  $p = .055$ ,  $\hat{\eta}_G^2 = .008$ . Post-hoc test revealed that the Good-Self condition

1366 ( $2.15 \pm 0.12$ ) was with greater  $d$  prime than Neutral condition ( $1.83 \pm 0.12$ ,  $t(34) = 3.36$ ,

1367  $p = 0.0031$ ), and Bad-self condition ( $1.87 \pm 0.12$ ),  $t(34) = 2.955$ ,  $p = 0.01$ . There was

1368 difference between neutral and bad condition,  $t(34) = -0.039$ ,  $p = 0.914$ . However, no effect

1369 of valence was found for the other-referential condition  $F(1.98, 114.61) = 1.75$ ,

1370  $MSE = 0.10$ ,  $p = .179$ ,  $\hat{\eta}_G^2 = .003$ .

1371 *Reaction time.*

1372 We found interaction between Matchness and Valence ( $F(1.94, 112.64) = 0.84$ ,

1373  $MSE = 465.35$ ,  $p = .432$ ,  $\hat{\eta}_G^2 = .000$ ) and then analyzed the matched trials and nonmatch

1374 trials separately, as in previous experiments.

1375 For the match trials, we found that the interaction between identity and valence,

1376  $F(1.90, 110.18) = 4.41$ ,  $MSE = 465.91$ ,  $p = .016$ ,  $\hat{\eta}_G^2 = .003$ , as well as the main effect of

1377 valence  $F(1.98, 114.82) = 0.94$ ,  $MSE = 606.30$ ,  $p = .392$ ,  $\hat{\eta}_G^2 = .001$ , but not the effect of

1378 identity  $F(1, 58) = 124.15$ ,  $MSE = 4,037.53$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .257$ . As for the  $d$  prime, we

1379 separated analyzed the self-referential and other-referential trials. For the Self-referential

1380 trials, we found the main effect of valence,  $F(1.97, 114.32) = 6.29$ ,  $MSE = 367.25$ ,

1381  $p = .003$ ,  $\hat{\eta}_G^2 = .006$ ; for the other-referential trials, the effect of valence is weaker,

1382  $F(1.95, 112.89) = 0.35$ ,  $MSE = 699.50$ ,  $p = .699$ ,  $\hat{\eta}_G^2 = .001$ . We then focused on the self

1383 conditions: the good-self condition ( $713 \pm 12$ ) is faster than neutral- ( $776 \pm 11.8$ ),  $t(34) =$

1384  $-7.396$ ,  $p < .0001$ , and bad-self ( $772 \pm 10.1$ ) conditions,  $t(34) = -5.66$ ,  $p < .0001$ . But

1385 there is not difference between neutral- and bad-self conditions,  $t(34) = 0.481$ ,  $p = 0.881$ .

1386 For the nonmatch trials, we didn't found any strong effect: identity,  $F(1, 58) = 0.16$ ,

1387  $MSE = 1,547.37$ ,  $p = .692$ ,  $\hat{\eta}_G^2 = .000$ , valence  $F(1.96, 113.52) = 0.68$ ,  $MSE = 390.26$ ,

1388  $p = .508$ ,  $\hat{\eta}_G^2 = .000$ , or interaction between the two  $F(1.90, 110.27) = 0.04$ ,

<sub>1389</sub>  $MSE = 585.80$ ,  $p = .953$ ,  $\hat{\eta}_G^2 = .000$ .

<sub>1390</sub> **BGLM.**

<sub>1391</sub> *Signal detection theory analysis of accuracy.*

<sub>1392</sub> We found that the  $d$  prime is greater when shapes were associated with good self  
<sub>1393</sub> condition than with neutral self or bad self, but shapes associated with bad self and neutral  
<sub>1394</sub> self didn't show differences. comparing the self vs other under three condition revealed that  
<sub>1395</sub> shapes associated with good self is greater than with good other, but with a weak evidence.  
<sub>1396</sub> In contrast, for both neutral and bad valence condition, shapes associated with other had  
<sub>1397</sub> greater  $d$  prime than with self.

<sub>1398</sub> *Reaction time.*

<sub>1399</sub> In reaction times, we found that same trends in the match trials as in the RT: while  
<sub>1400</sub> the shapes associated with good self was greater than with good other (log mean diff =  
<sub>1401</sub> -0.02858, 95%HPD[-0.070898, 0.0154]), the direction is reversed for neutral and negative  
<sub>1402</sub> condition. see Figure 12

<sub>1403</sub> **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et  
<sub>1404</sub> al., 2020). We estimated separate drift rate ( $v$ ), non-decision time ( $T_0$ ), and boundary  
<sub>1405</sub> separation ( $a$ ) for each condition. We found that the shapes tagged with good person has  
<sub>1406</sub> higher drift rate and higher boundary separation than shapes tagged with both neutral and  
<sub>1407</sub> bad person. Also, the shapes tagged with neutral person has a higher drift rate than  
<sub>1408</sub> shapes tagged with bad person, but not for the boundary separation. Finally, we found  
<sub>1409</sub> that shapes tagged with bad person had longer non-decision time (see figure 13)).

<sub>1410</sub> **Experiment 4b: Morality as task-irrelevant variable**

<sub>1411</sub> In study 4b, we changed the role of valence and identity in task. In this experiment,  
<sub>1412</sub> participants learn the association between moral valence and the made perceptual match  
<sub>1413</sub> judgments to associations between different moral valence and shapes as in study 1-3.

<sup>1414</sup> Different from experiment 1 ~ 3, we made put the labels of “self/other” in the shapes so  
<sup>1415</sup> that identity served as an task irrelevant variable. As in experiment 4b, we also  
<sup>1416</sup> hypothesized that the intrinsic binding between morally good and self will enhance the  
<sup>1417</sup> performance of good self condition, even identity is irrelevant to the task.

<sup>1418</sup> **Method.**

<sup>1419</sup> ***Participants.***

<sup>1420</sup> 53 participants (39 female, age =  $20.57 \pm 1.81$ ) participated the current study, 34 of  
<sup>1421</sup> them were from Tsinghua University in 2015, 19 were from Wenzhou University  
<sup>1422</sup> participated in 2017. All participants were right-handed, and all had normal or  
<sup>1423</sup> corrected-to-normal vision. Informed consent was obtained from all participants prior to  
<sup>1424</sup> the experiment according to procedures approved by a local ethics committee. The data  
<sup>1425</sup> from 8 participants from Wenzhou site were excluded from analysis because their accuracy  
<sup>1426</sup> was close to chance ( $< 0.6$ ). The results for the remaining 45 participants (33 female, age  
<sup>1427</sup> =  $20.78 \pm 1.76$ ) were analyzed and reported.

<sup>1428</sup> ***Design.***

<sup>1429</sup> As in Experiment 3, a  $2 \times 3 \times 2$  within-subject design was used. The first variable was  
<sup>1430</sup> self-relevance (self and stranger associations); the second variable was moral valence (good,  
<sup>1431</sup> neutral and bad associations); the third variable was the matching between shape and label  
<sup>1432</sup> (matching vs. non-match for the personal association). However, in this the task,  
<sup>1433</sup> participants only learn the association between two geometric shapes and two labels (self  
<sup>1434</sup> and other), i.e., only self-relevance were related to the task. The moral valence  
<sup>1435</sup> manipulation was achieved by embedding the personal label of the labels in the geometric  
<sup>1436</sup> shapes, see below. For simplicity, the trials where shapes where paired with self and with a  
<sup>1437</sup> word of “good person” inside were shorted as good-self condition, similarly, the trials where  
<sup>1438</sup> shapes paired with the self and with a word of “bad person” inside were shorted as bad-self  
<sup>1439</sup> condition. Hence, we also have six conditions: good-self, neutral-self, bad-self, good-other,

1440 neutral-other, and bad-other.

1441 ***Stimuli.***

1442 2 shapes were included (circle, square) and each appeared above a central fixation  
1443 cross with the personal label appearing below. However, the shapes were not empty but  
1444 with a two-Chinese-character word in the middle, the word was one of three labels with  
1445 different moral valence: “good person”, “bad person” and “neutral person”. Before the  
1446 experiment, participants learned the self/other association, and were informed to only  
1447 response to the association between shapes’ configures and the labels below the fixation, but  
1448 ignore the words within shapes. Besides the behavioral experiments, participants from  
1449 Tsinghua community also finished questionnaires as Experiments 3, and participants from  
1450 Wenzhou community finished a series of questionnaire as the other experiment finished in  
1451 Wenzhou.

1452 ***Procedure.***

1453 The procedure was similar to Experiment 1. There were 6 blocks of trial, each with  
1454 120 trials for 2017 data. Due to procedure error, the data collected in 2015 in Tsinghua  
1455 community only have 60 trials for each block, i.e., 30 trials per condition.

1456 As in study 3a, before each task, the instruction showed the meaning of each label to  
1457 participants. The self-matching task and other-matching task were randomized between  
1458 participants. Each participant finished 6 blocks, each have 120 trials.

1459 ***Data Analysis.***

1460 Same as experiment 3a.

1461 **Results.**

1462 ***NHST.***

1463 Figure 14 shows  $d'$  prime and reaction times of experiment 3a. Less than 5% correct  
1464 trials with less than 200ms reaction times were excluded.

<sub>1465</sub> *d prime.*

<sub>1466</sub> There was no evidence for the main effect of valence,  $F(1.59, 69.94) = 2.34$ ,

<sub>1467</sub>  $MSE = 0.48$ ,  $p = .115$ ,  $\hat{\eta}_G^2 = .010$ , but we found a main effect of self-relevance,

<sub>1468</sub>  $F(1, 44) = 0.00$ ,  $MSE = 0.08$ ,  $p = .994$ ,  $\hat{\eta}_G^2 = .000$ , as well as the interaction,

<sub>1469</sub>  $F(1.96, 86.41) = 0.53$ ,  $MSE = 0.10$ ,  $p = .585$ ,  $\hat{\eta}_G^2 = .001$ .

<sub>1470</sub> We then conducted separated ANOVA for self-referential and other-referential trials.

<sub>1471</sub> The valence effect was shown for the self-referential conditions,  $F(1.75, 76.86) = 3.08$ ,

<sub>1472</sub>  $MSE = 0.25$ ,  $p = .058$ ,  $\hat{\eta}_G^2 = .017$ . Post-hoc test revealed that the Good-Self condition

<sub>1473</sub> ( $2.15 \pm 0.12$ ) was with greater *d* prime than Neutral condition ( $1.83 \pm 0.12$ ,  $t(34) = 3.36$ ,

<sub>1474</sub>  $p = 0.0031$ ), and Bad-self condition ( $1.87 \pm 0.12$ ),  $t(34) = 2.955$ ,  $p = 0.01$ . There was

<sub>1475</sub> difference between neutral and bad condition,  $t(34) = -0.039$ ,  $p = 0.914$ . However, no effect

<sub>1476</sub> of valence was found for the other-referential condition  $F(1.63, 71.50) = 1.07$ ,  $MSE = 0.33$ ,

<sub>1477</sub>  $p = .336$ ,  $\hat{\eta}_G^2 = .006$ .

<sub>1478</sub> *Reaction time.*

<sub>1479</sub> We found interaction between Matchness and Valence ( $F(1.87, 82.50) = 18.58$ ,

<sub>1480</sub>  $MSE = 1,291.12$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .023$ ) and then analyzed the matched trials and

<sub>1481</sub> nonmatch trials separately, as in previous experiments.

<sub>1482</sub> For the match trials, we found that the interaction between identity and valence,

<sub>1483</sub>  $F(1.86, 81.84) = 5.22$ ,  $MSE = 308.30$ ,  $p = .009$ ,  $\hat{\eta}_G^2 = .003$ , as well as the main effect of

<sub>1484</sub> valence  $F(1.80, 79.37) = 11.04$ ,  $MSE = 2,937.54$ ,  $p < .001$ ,  $\hat{\eta}_G^2 = .059$ , but not the effect of

<sub>1485</sub> identity  $F(1, 44) = 0.23$ ,  $MSE = 263.26$ ,  $p = .632$ ,  $\hat{\eta}_G^2 = .000$ . As for the *d* prime, we

<sub>1486</sub> separated analyzed the self-referential and other-referential trials. For the Self-referential

<sub>1487</sub> trials, we found the main effect of valence,  $F(1.74, 76.48) = 13.69$ ,  $MSE = 1,732.08$ ,

<sub>1488</sub>  $p < .001$ ,  $\hat{\eta}_G^2 = .079$ ; for the other-referential trials, the effect of valence is weaker,

<sub>1489</sub>  $F(1.87, 82.44) = 7.09$ ,  $MSE = 1,527.43$ ,  $p = .002$ ,  $\hat{\eta}_G^2 = .043$ . We then focused on the self

<sub>1490</sub> conditions: the good-self condition ( $713 \pm 12$ ) is faster than neutral- ( $776 \pm 11.8$ ),  $t(34) =$

<sup>1491</sup> -7.396,  $p < .0001$ , and bad-self ( $772 \pm 10.1$ ) conditions,  $t(34) = -5.66$ ,  $p < .0001$ . But  
<sup>1492</sup> there is not difference between neutral- and bad-self conditions,  $t(34) = 0.481$ ,  $p = 0.881$ .

<sup>1493</sup> For the nonmatch trials, we didn't found any strong effect: identity,  $F(1, 44) = 1.96$ ,  
<sup>1494</sup>  $MSE = 319.47$ ,  $p = .169$ ,  $\hat{\eta}_G^2 = .001$ , valence  $F(1.69, 74.54) = 6.59$ ,  $MSE = 886.19$ ,  
<sup>1495</sup>  $p = .004$ ,  $\hat{\eta}_G^2 = .010$ , or interaction between the two  $F(1.88, 82.57) = 0.31$ ,  $MSE = 316.96$ ,  
<sup>1496</sup>  $p = .718$ ,  $\hat{\eta}_G^2 = .000$ .

<sup>1497</sup> **BGLM.**

<sup>1498</sup> *Signal detection theory analysis of accuracy.*

<sup>1499</sup> We found that the  $d$  prime is greater when shapes were associated with good self  
<sup>1500</sup> condition than with neutral self or bad self, but shapes associated with bad self and neutral  
<sup>1501</sup> self didn't show differences. comparing the self vs other under three condition revealed that  
<sup>1502</sup> shapes associated with good self is greater than with good other, but with a weak evidence.  
<sup>1503</sup> In contrast, for both neutral and bad valence condition, shapes associated with other had  
<sup>1504</sup> greater  $d$  prime than with self.

<sup>1505</sup> *Reaction time.*

<sup>1506</sup> In reaction times, we found that same trends in the match trials as in the RT: while  
<sup>1507</sup> the shapes associated with good self was greater than with good other (log mean diff =  
<sup>1508</sup>  $-0.02858$ , 95%HPD[-0.070898, 0.0154]), the direction is reversed for neutral and negative  
<sup>1509</sup> condition. see Figure 15

<sup>1510</sup> **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et  
<sup>1511</sup> al., 2020). We estimated separate drift rate ( $v$ ), non-decision time ( $T_0$ ), and boundary  
<sup>1512</sup> separation ( $a$ ) for each condition. We found that the shapes tagged with good person has  
<sup>1513</sup> higher drift rate and higher boundary separation than shapes tagged with both neutral and  
<sup>1514</sup> bad person. Also, the shapes tagged with neutral person has a higher drift rate than  
<sup>1515</sup> shapes tagged with bad person, but not for the boundary separation. Finally, we found  
<sup>1516</sup> that shapes tagged with bad person had longer non-decision time (see figure 16)).

1517

## Results

1518 **Effect of moral valence**

1519 In this part, we synthesized results from experiment 1a, 1b, 1c, 2, 5 and 6a. Data  
1520 from 192 participants were included in these analyses. We found differences between  
1521 positive and negative conditions on RT was Cohen's  $d = -0.58 \pm 0.06$ , 95% CI [-0.70 -0.47];  
1522 on  $d'$  was Cohen's  $d = 0.24 \pm 0.05$ , 95% CI [0.15 0.34]. The effect was also observed  
1523 between positive and neutral condition, RT: Cohen's  $d = -0.44 \pm 0.10$ , 95% CI [-0.63  
1524 -0.25];  $d'$ : Cohen's  $d = 0.31 \pm 0.07$ , 95% CI [0.16 0.45]. And the difference between neutral  
1525 and bad conditions are not significant, RT: Cohen's  $d = 0.15 \pm 0.07$ , 95% CI [0.00 0.30];  
1526  $d'$ : Cohen's  $d = 0.07 \pm 0.07$ , 95% CI [-0.08 0.21]. See Figure 17 left panel.

1527 **Interaction between valence and self-reference**

1528 In this part, we combined the experiments that explicitly manipulated the  
1529 self-reference and valence, which includes 3a, 3b, 6b, 7a, and 7b. For the positive versus  
1530 negative contrast, data were from five experiments with 178 participants; for positive  
1531 versus neutral and neutral versus negative contrasts, data were from three experiments ( (

1532 3a, 3b, and 6b) with 108 participants.

1533 In most of these experiments, the interaction between self-reference and valence was  
1534 significant (see results of each experiment in supplementary materials). In the  
1535 mini-meta-analysis, we analyzed the valence effect for self-referential condition and  
1536 other-referential condition separately.

1537 For the self-referential condition, we found the same pattern as in the first part of  
1538 results. That is we found significant differences between positive and neutral as well as  
1539 positive and negative, but not neutral and negative. The effect size of RT between positive  
1540 and negative is Cohen's  $d = -0.89 \pm 0.12$ , 95% CI [-1.11 -0.66]; on  $d'$  was Cohen's  $d = 0.61$

1541  $\pm 0.09$ , 95% CI [0.44 0.78]. The effect was also observed between positive and neutral  
1542 condition, RT: Cohen's  $d = -0.76 \pm 0.13$ , 95% CI [-1.01 -0.50];  $d'$ : Cohen's  $d = 0.69 \pm$   
1543 0.14, 95% CI [0.42 0.96]. And the difference between neutral and bad conditions are not  
1544 significant, RT: Cohen's  $d = 0.03 \pm 0.13$ , 95% CI [-0.22 0.29];  $d'$ : Cohen's  $d = 0.08 \pm 0.08$ ,  
1545 95% CI [-0.07 0.24]. See Figure 17 the middle panel.

1546 For the other-referential condition, we found that only the difference between positive  
1547 and negative on RT was significant, all the other conditions were not. The effect size of RT  
1548 between positive and negative is Cohen's  $d = -0.28 \pm 0.05$ , 95% CI [-0.38 -0.17]; on  $d'$  was  
1549 Cohen's  $d = -0.02 \pm 0.08$ , 95% CI [-0.17 0.13]. The effect was not observed between  
1550 positive and neutral condition, RT: Cohen's  $d = -0.12 \pm 0.10$ , 95% CI [-0.31 0.06];  $d'$ :  
1551 Cohen's  $d = 0.01 \pm 0.08$ , 95% CI [-0.16 0.17]. And the difference between neutral and bad  
1552 conditions are not significant, RT: Cohen's  $d = 0.14 \pm 0.09$ , 95% CI [-0.03 0.31];  $d'$ :  
1553 Cohen's  $d = 0.05 \pm 0.07$ , 95% CI [-0.08 0.18]. See Figure 17 right panel.

#### 1554 Generalizability of the valence effect

1555 In this part, we reported the results from experiment 4 in which either moral valence  
1556 or self-reference were manipulated as task-irrelevant stimuli.

1557 For experiment 4a, when self-reference was the target and moral valence was  
1558 task-irrelevant, we found that only under the implicit self-referential condition, i.e., when  
1559 the moral words were presented as task irrelevant stimuli, there was the main effect of  
1560 valence and interaction between valence and reference for both  $d$  prime and RT (See  
1561 supplementary results for the detailed statistics). For  $d$  prime, we found good-self  
1562 condition ( $2.55 \pm 0.86$ ) had higher  $d$  prime than bad-self condition ( $2.38 \pm 0.80$ ); good self  
1563 condition was also higher than neutral self ( $2.45 \pm 0.78$ ) but there was not statistically  
1564 significant, while the neutral-self condition was higher than bad self condition and not  
1565 significant neither. For reaction times, good-self condition ( $654.26 \pm 67.09$ ) were faster

1566 relative to bad-self condition ( $665.64 \pm 64.59$ ), and over neutral-self condition ( $664.26 \pm$   
1567  $64.71$ ). The difference between neutral-self and bad-self conditions were not significant.  
1568 However, for the other-referential condition, there was no significant differences between  
1569 different valence conditions. See Figure 18.

1570 For experiment 4b, when valence was the target and the identity was task-irrelevant,  
1571 we found a strong valence effect (see supplementary results and Figure 19, Figure 20).

1572 In this experiment, the advantage of good-self condition can only be disentangled by  
1573 comparing the self-referential and other-referential conditions. Therefore, we calculated the  
1574 differences between the valence effect under self-referential and other referential conditions  
1575 and used the weighted variance as the variance of this differences. We found this  
1576 modulation effect on RT. The valence effect of RT was stronger in self-referential than  
1577 other-referential for the Good vs. Neutral condition ( $-0.33 \pm 0.01$ ), and to a less extent the  
1578 Good vs. Bad condition ( $-0.17 \pm 0.01$ ). While the size of the other effect's CI included  
1579 zero, suggestion those effects didn't differ from zero. See Figure 21.

## 1580 Specificity of valence effect

1581 In this part, we analyzed the results from experiment 5, which included positive,  
1582 neutral, and negative valence from four different domains: morality, emotion, aesthetics of  
1583 human, and aesthetics of scene. We found interaction between valence and domain for both  
1584 *d* prime and RT (match trials). A common pattern appeared in all four domains: each  
1585 domain showed a binary results instead of gradient on both *d* prime and RT. For morality,  
1586 aesthetics of human, and aesthetics of scene, there was a positivity effect where the positive  
1587 conditions had advantages over both neutral (greater *d* prime and faster RT), while neutral  
1588 and negative conditions didn't differ from each other. But for the emotional stimuli, there  
1589 was a reversed negativity effect: positive and neutral conditions were not significantly  
1590 different from each other but both had advantage over negative conditions. See

supplementary materials for detailed statistics. Also note that the effect size in moral domain is smaller than the aesthetic domains (beauty of people and beauty of scene). See Figure 22.

#### Self-reported personal distance

See Figure 23.

#### Correlation analyses

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

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

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

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

First, boundary separation ( $a$ ) of moral good condition was correlated with both Self-Bad distance ( $r = 0.198$ , 95% CI [],  $p = 0.0063$ ) and Neutral-Bad distance

1615 ( $r = 0.1571$ , 95% CI [],  $p = 0.031$ ). At the same time, the non-decision time is negatively  
1616 correlated with Self-Bad distance ( $r = 0.169$ , 95% CI [],  $p = 0.0197$ ). See Figure 24.

1617 Second, we found the boundary separation of neutral condition is positively  
1618 correlated with the personal distance between self and good distance ( $r = 0.189$ , 95% CI [],  
1619  $p = 0.036$ ), but negatively correlated with self-neutral distance( $r = -0.183$ , 95% CI [],  
1620  $p = 0.042$ ). Also, the drift rate of the neutral condition is positively correlated with the  
1621 Self-Bad distance ( $r = 0.177$ , 95% CI [],  $p = 0.048$ ).a. See figure 25

1622 We also explored the correlation between behavioral data and questionnaire scores  
1623 separately for experiments with and without self-referential, however, the sample size is  
1624 very low for some conditions.

1625 **Discussion**

1626 **References**

- 1627 Abele, A. E., Ellemers, N., Fiske, S. T., Koch, A., & Yzerbyt, V. (2020). Navigating the  
1628 social world: Toward an integrated framework for evaluating self, individuals, and  
1629 groups. *Psychological Review*. <https://doi.org/10.1037/rev0000262>
- 1630 Amodio, D. M. (2019). Social cognition 2.0: An interactive memory systems account.  
1631 *Trends in Cognitive Sciences*, 23(1), 21–33.  
1632 <https://doi.org/10.1016/j.tics.2018.10.002>
- 1633 Anderson, E., Siegel, E. H., Bliss-Moreau, E., & Barrett, L. F. (2011). The visual impact  
1634 of gossip. *Science*, 332(6036), 1446–1448. <https://doi.org/10.1126/science.1201574>
- 1635 Brainard, D. H. (1997). The psychophysics toolbox. *Spatial Vision*, 10(4), 433–436.  
1636 Journal Article.
- 1637 Bürkner, P.-C. (2017). Brms: An r package for bayesian multilevel models using stan.  
1638 *Journal of Statistical Software*; Vol 1, Issue 1 (2017). Journal Article. Retrieved

- 1639 from  
1640 <https://www.jstatsoft.org/v080/i01> <http://dx.doi.org/10.18637/jss.v080.i01>
- 1641 Carlson, R. W., Maréchal, M. A., Oud, B., Fehr, E., & Crockett, M. J. (2020). Motivated  
1642 misremembering of selfish decisions. *Nature Communications*, 11(1), 2100.  
1643 <https://doi.org/10.1038/s41467-020-15602-4>
- 1644 Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M., ...  
1645 Riddell, A. (2017). Stan: A probabilistic programming language. *Journal of  
1646 Statistical Software*, 76(1). Journal Article. <https://doi.org/10.18637/jss.v076.i01>
- 1647 Cooper, H., Hedges, L. V., & Valentine, J. C. (2009). *The handbook of research synthesis  
1648 and meta-analysis* (2nd ed.). Book, New York: Sage.
- 1649 Dang, J., King, K. M., & Inzlicht, M. (2020). Why are self-report and behavioral measures  
1650 weakly correlated? *Trends in Cognitive Sciences*, 24(4), 267–269.  
1651 <https://doi.org/10.1016/j.tics.2020.01.007>
- 1652 DeCarlo, L. T. (1998). Signal detection theory and generalized linear models. *Psychological  
1653 Methods*, 3(2), 186–205. Journal Article. <https://doi.org/10.1037/1082-989X.3.2.186>
- 1654 Dzhelyova, M., Perrett, D. I., & Jentzsch, I. (2012). Temporal dynamics of trustworthiness  
1655 perception. *Brain Research*, 1435, 81–90.  
1656 <https://doi.org/10.1016/j.brainres.2011.11.043>
- 1657 Ellemers, N., Toorn, J. van der, Paunov, Y., & Leeuwen, T. van. (2019). The psychology of  
1658 morality: A review and analysis of empirical studies published from 1940 through  
1659 2017. *Personality and Social Psychology Review*, 23(4), 332–366.  
1660 <https://doi.org/10.1177/1088868318811759>
- 1661 Enock, F. E., Hewstone, M. R. C., Lockwood, P. L., & Sui, J. (2020). Overlap in  
1662 processing advantages for minimal ingroups and the self. *Scientific Reports*, 10(1),  
1663 18933. <https://doi.org/10.1038/s41598-020-76001-9>

- 1664 Enock, F., Sui, J., Hewstone, M., & Humphreys, G. W. (2018). Self and team prioritisation  
1665 effects in perceptual matching: Evidence for a shared representation. *Acta  
1666 Psychologica*, 182, 107–118. <https://doi.org/10.1016/j.actpsy.2017.11.011>
- 1667 Faul, F., Erdfelder, E., Buchner, A., & Lang, A.-G. (2009). Statistical power analyses using  
1668 g\*Power 3.1: Tests for correlation and regression analyses. *Behavior Research  
1669 Methods*, 41(4), 1149–1160. <https://doi.org/10.3758/BRM.41.4.1149>
- 1670 Firestone, C., & Scholl, B. J. (2015). Enhanced visual awareness for morality and pajamas?  
1671 Perception vs. Memory in “top-down” effects. *Cognition*, 136, 409–416.  
1672 <https://doi.org/10.1016/j.cognition.2014.10.014>
- 1673 Freeman, J. B., & Ambady, N. (2011). A dynamic interactive theory of person construal.  
1674 *Psychological Review*, 118(2), 247–279. <https://doi.org/10.1037/a0022327>
- 1675 Gantman, A. P., & Van Bavel, J. J. (2014). The moral pop-out effect: Enhanced  
1676 perceptual awareness of morally relevant stimuli. *Cognition*, 132(1), 22–29.  
1677 <https://doi.org/10.1016/j.cognition.2014.02.007>
- 1678 Goh, J. X., Hall, J. A., & Rosenthal, R. (2016). Mini meta-analysis of your own studies:  
1679 Some arguments on why and a primer on how. *Social and Personality Psychology  
1680 Compass*, 10(10), 535–549. Journal Article. <https://doi.org/10.1111/spc3.12267>
- 1681 Goodwin, G. P. (2015). Moral character in person perception. *Current Directions in  
1682 Psychological Science*, 24(1), 38–44. <https://doi.org/10.1177/0963721414550709>
- 1683 Goodwin, G. P., Piazza, J., & Rozin, P. (2014). Moral character predominates in person  
1684 perception and evaluation. *Journal of Personality and Social Psychology*, 106(1),  
1685 148–168. <https://doi.org/10.1037/a0034726>
- 1686 Henrich, J., Heine, S. J., & Norenzayan, A. (2010). The weirdest people in the world?  
1687 *Behavioral and Brain Sciences*, 33(2), 61–83.  
1688 <https://doi.org/10.1017/S0140525X0999152X>

- 1689 Hofmann, W., Wisneski, D. C., Brandt, M. J., & Skitka, L. J. (2014). Morality in everyday  
1690 life. *Science*, 345(6202), 1340–1343. <https://doi.org/10.1126/science.1251560>
- 1691 Hu, C.-P., Lan, Y., Macrae, C. N., & Sui, J. (2020). Good me bad me: Does valence  
1692 influence self-prioritization during perceptual decision-making? *Collabra: Psychology*,  
1693 6(1), 20. Journal Article. <https://doi.org/10.1525/collabra.301>
- 1694 Juechems, K., & Summerfield, C. (2019). Where does value come from? *Trends in  
1695 Cognitive Sciences*, 23(10), 836–850. <https://doi.org/10.1016/j.tics.2019.07.012>
- 1696 Kirby, K. N., & Gerlanc, D. (2013). BootES: An r package for bootstrap confidence  
1697 intervals on effect sizes. *Behavior Research Methods*, 45(4), 905–927.  
1698 <https://doi.org/10.3758/s13428-013-0330-5>
- 1699 Klein, N., & Epley, N. (2016). Maybe holier, but definitely less evil, than you: Bounded  
1700 self-righteousness in social judgment. *Journal of Personality and Social Psychology*,  
1701 110(5), 660–674. <https://doi.org/10.1037/pspa0000050>
- 1702 Liu, Q., Wang, F., Yan, W., Peng, K., Sui, J., & Hu, C.-P. (2020). Questionnaire data from  
1703 the revision of a chinese version of free will and determinism plus scale. *Journal of  
1704 Open Psychology Data*, 8(1), 1. Journal Article. <https://doi.org/10.5334/jopd.49/>
- 1705 Matzke, D., & Wagenmakers, E.-J. (2009). Psychological interpretation of the ex-gaussian  
1706 and shifted wald parameters: A diffusion model analysis. *Psychonomic Bulletin &  
1707 Review*, 16(5), 798–817. <https://doi.org/10.3758/PBR.16.5.798>
- 1708 McHugh, C., McGann, M., Igou, E. R., & Kinsella, E. (2019). *Moral judgment as  
1709 categorization (MJAC)*. PsyArXiv. <https://doi.org/10.31234/osf.io/72dzp>
- 1710 McNeish, D., & Wolf, M. G. (2020). Thinking twice about sum scores. *Behavior Research  
1711 Methods*. <https://doi.org/10.3758/s13428-020-01398-0>
- 1712 Monin, B., & Jordan, A. H. (2009). The dynamic moral self: A social psychological  
1713 perspective. In *Personality, identity, and character: Explorations in moral*

- 1714 psychology (pp. 341–354). New York, NY, US: Cambridge University Press.
- 1715 <https://doi.org/10.1017/CBO9780511627125.016>
- 1716 Pelli, D. G. (1997). The videotoolbox software for visual psychophysics: Transforming  
1717 numbers into movies. *Spatial Vision*, 10(4), 437–442. Journal Article.
- 1718 Reicher, S., & Hopkins, N. (2016). Perception, action, and the social dynamics of the  
1719 variable self. *Psychological Inquiry*, 27(4), 341–347.  
1720 <https://doi.org/10.1080/1047840X.2016.1217584>
- 1721 Rouder, J. N., & Lu, J. (2005). An introduction to bayesian hierarchical models with an  
1722 application in the theory of signal detection. *Psychonomic Bulletin & Review*,  
1723 12(4), 573–604. Journal Article. <https://doi.org/10.3758/bf03196750>
- 1724 Rousselet, G. A., & Wilcox, R. R. (2019). Reaction times and other skewed distributions:  
1725 Problems with the mean and the median. *Meta-Psychology*. preprint.  
1726 <https://doi.org/10.1101/383935>
- 1727 Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2013). Life after p-hacking. Conference  
1728 Proceedings. <https://doi.org/10.2139/ssrn.2205186>
- 1729 Stanislaw, H., & Todorov, N. (1999). Calculation of signal detection theory measures.  
1730 *Behavior Research Methods, Instruments, & Computers*, 31(1), 137–149. Journal  
1731 Article. <https://doi.org/10.3758/BF03207704>
- 1732 Stanley, M. L., & De Brigard, F. (2019). Moral memories and the belief in the good self.  
1733 *Current Directions in Psychological Science*, 28(4), 387–391.  
1734 <https://doi.org/10.1177/0963721419847990>
- 1735 Stein, T., Grubb, C., Bertrand, M., Suh, S. M., & Verosky, S. C. (2017). No impact of  
1736 affective person knowledge on visual awareness: Evidence from binocular rivalry and  
1737 continuous flash suppression. *Emotion*, 17(8), 1199–1207.  
1738 <https://doi.org/10.1037/emo0000305>

- 1739 Stolier, R. M., & Freeman, J. B. (2016). Functional and temporal considerations for  
1740 top-down influences in social perception. *Psychological Inquiry*, 27(4), 352–357.  
1741 <https://doi.org/10.1080/1047840X.2016.1216034>
- 1742 Strohminger, N., Knobe, J., & Newman, G. (2017). The true self: A psychological concept  
1743 distinct from the self: *Perspectives on Psychological Science*.  
1744 <https://doi.org/10.1177/1745691616689495>
- 1745 Sui, J., He, X., & Humphreys, G. W. (2012). Perceptual effects of social salience: Evidence  
1746 from self-prioritization effects on perceptual matching. *Journal of Experimental  
1747 Psychology: Human Perception and Performance*, 38(5), 1105–1117. Journal  
1748 Article. <https://doi.org/10.1037/a0029792>
- 1749 Tappin, B. M., & McKay, R. T. (2017). The illusion of moral superiority. *Social  
1750 Psychological and Personality Science*, 8(6), 623–631.  
1751 <https://doi.org/10.1177/1948550616673878>
- 1752 Turner, J. C., Hogg, M. A., Oakes, P. J., Reicher, S. D., & Wetherell, M. S. (1987).  
1753 *Rediscovering the social group: A self-categorization theory*. Cambridge, MA, US:  
1754 Basil Blackwell.
- 1755 Turner, J. C., Oakes, P. J., Haslam, S. A., & McGarty, C. (1994). Self and collective:  
1756 Cognition and social context. *Personality and Social Psychology Bulletin*, 20(5),  
1757 454–463. <https://doi.org/10.1177/0146167294205002>
- 1758 Uhlmann, E. L., Pizarro, D. A., & Diermeier, D. (2015). A person-centered approach to  
1759 moral judgment: *Perspectives on Psychological Science*.  
1760 <https://doi.org/10.1177/1745691614556679>
- 1761 Verosky, S. C., & Todorov, A. (2010). Differential neural responses to faces physically  
1762 similar to the self as a function of their valence. *NeuroImage*, 49(2), 1690–1698.  
1763 <https://doi.org/10.1016/j.neuroimage.2009.10.017>

- 1764 Waytz, A., Dungan, J., & Young, L. (2013). The whistleblower's dilemma and the  
1765 fairness–loyalty tradeoff. *Journal of Experimental Social Psychology*, 49(6),  
1766 1027–1033. <https://doi.org/10.1016/j.jesp.2013.07.002>
- 1767 Wiecki, T. V., Sofer, I., & Frank, M. J. (2013). HDDM: Hierarchical bayesian estimation of  
1768 the drift-diffusion model in python. *Frontiers in Neuroinformatics*, 7.  
1769 <https://doi.org/10.3389/fninf.2013.00014>
- 1770 Willis, J., & Todorov, A. (2006). First impressions: Making up your mind after a 100-ms  
1771 exposure to a face. *Psychological Science*, 17(7), 592–598.  
1772 <https://doi.org/10.1111/j.1467-9280.2006.01750.x>
- 1773 Xiao, Y. J., Coppin, G., & Bavel, J. J. V. (2016). Perceiving the world through  
1774 group-colored glasses: A perceptual model of intergroup relations. *Psychological  
1775 Inquiry*, 27(4), 255–274. <https://doi.org/10.1080/1047840X.2016.1199221>

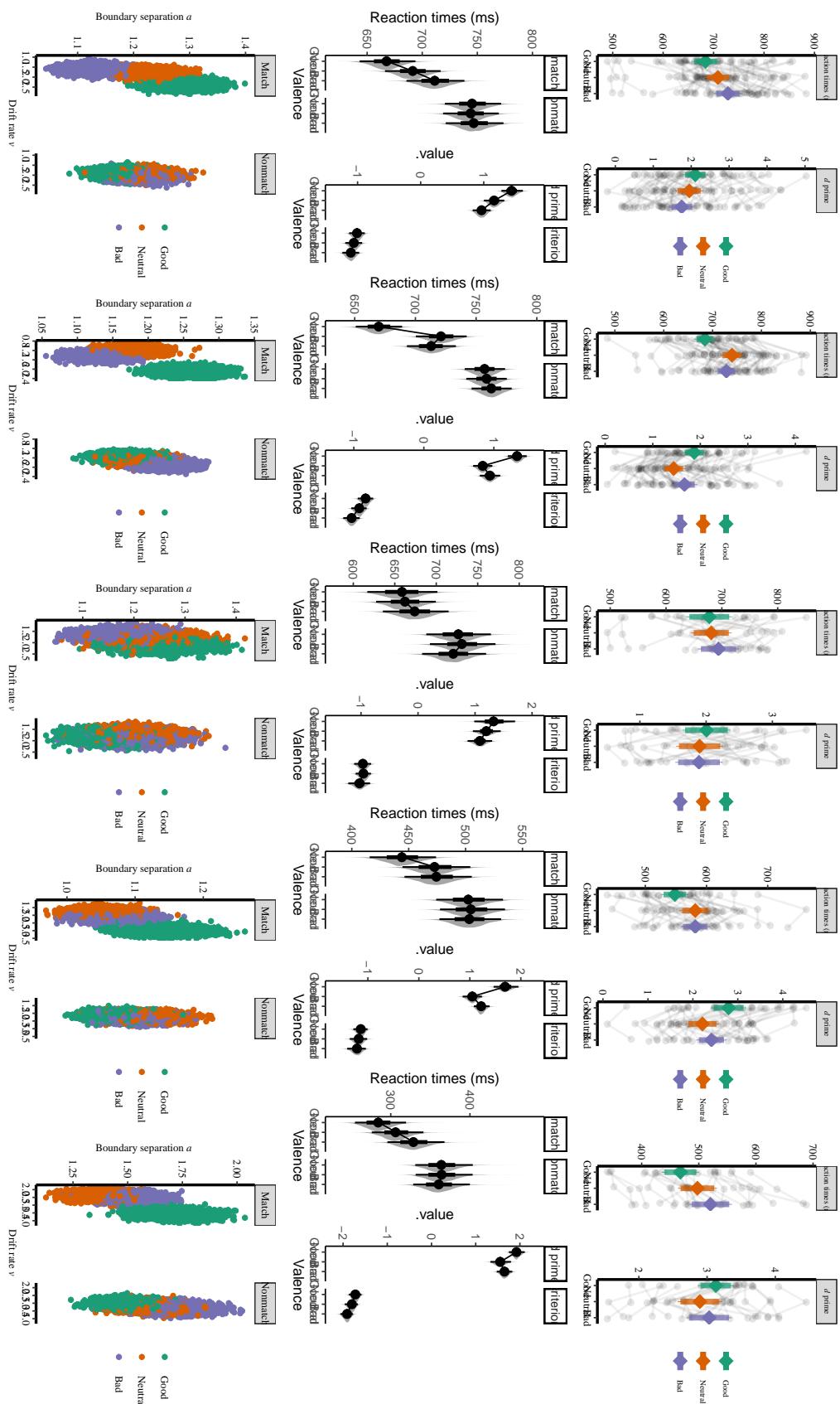


Figure 1. Results for part 1.

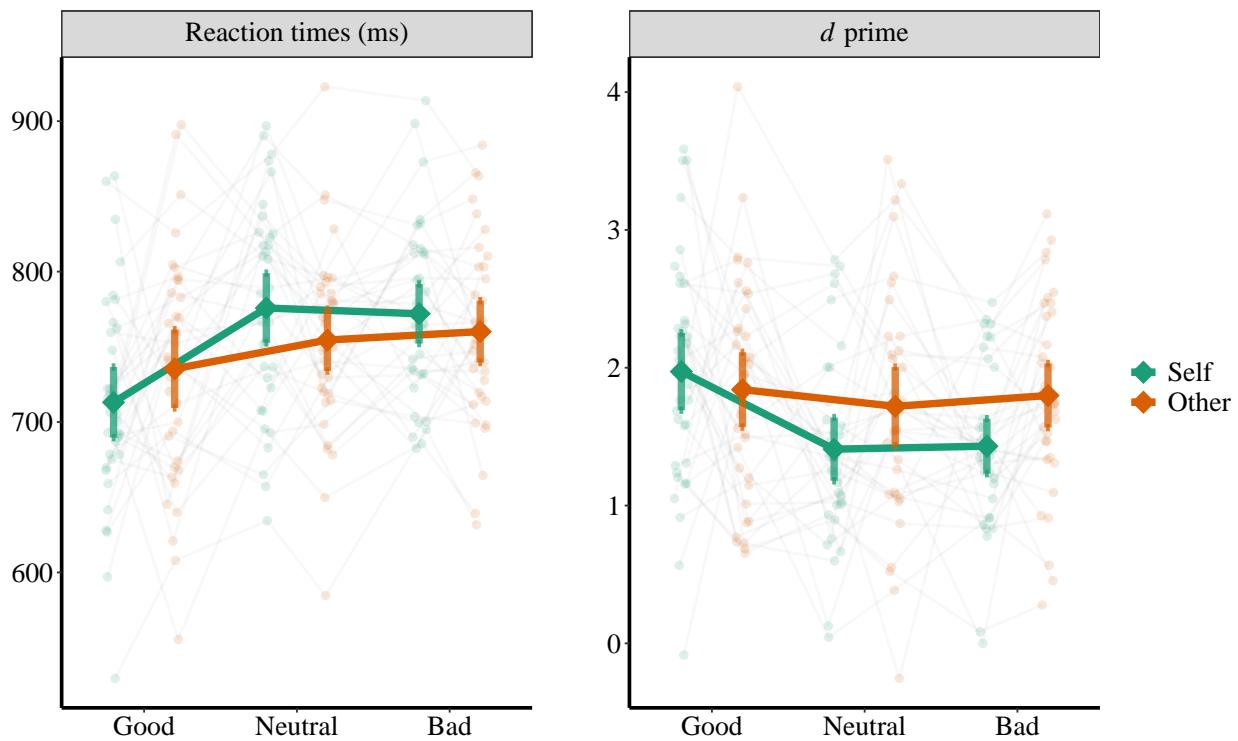


Figure 2. RT and  $d$  prime of Experiment 3a.

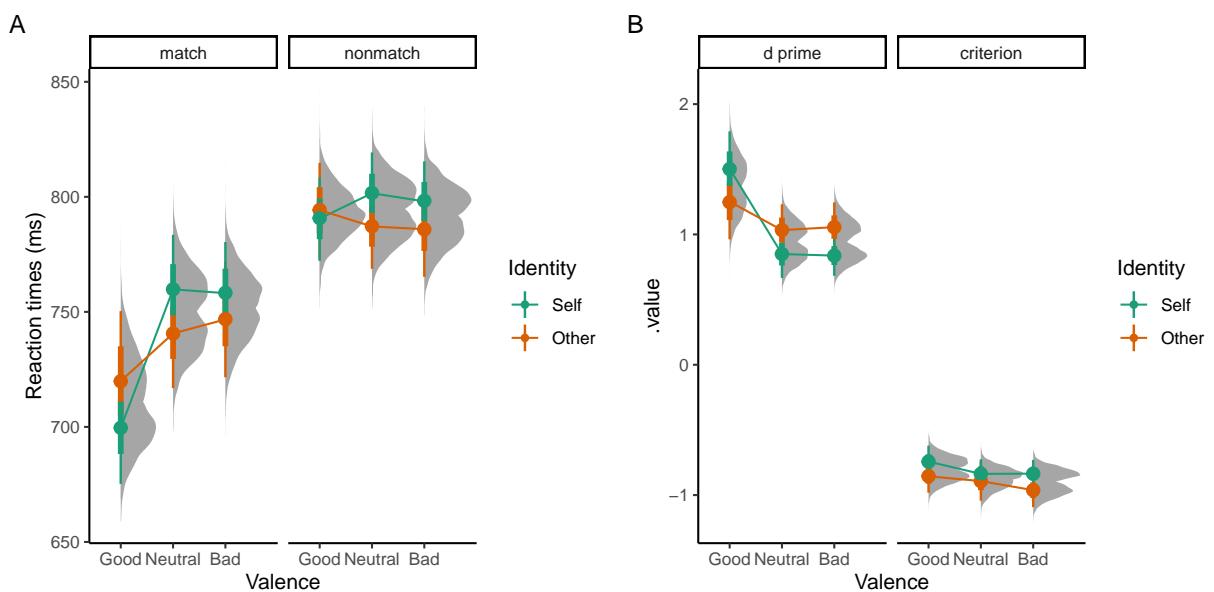


Figure 3. Exp3a: Results of Bayesian GLM analysis.

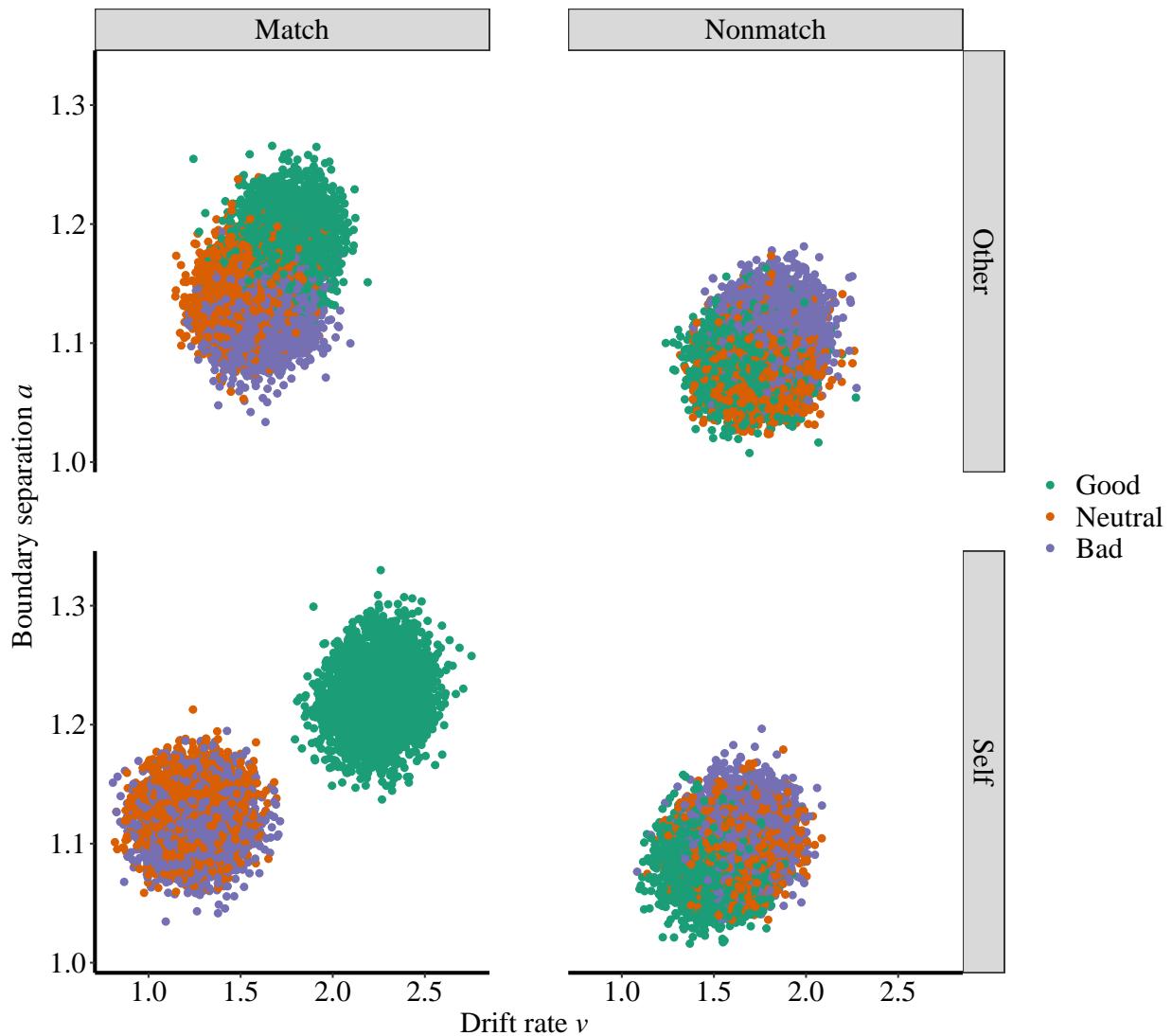


Figure 4. Exp3a: Results of HDDM.

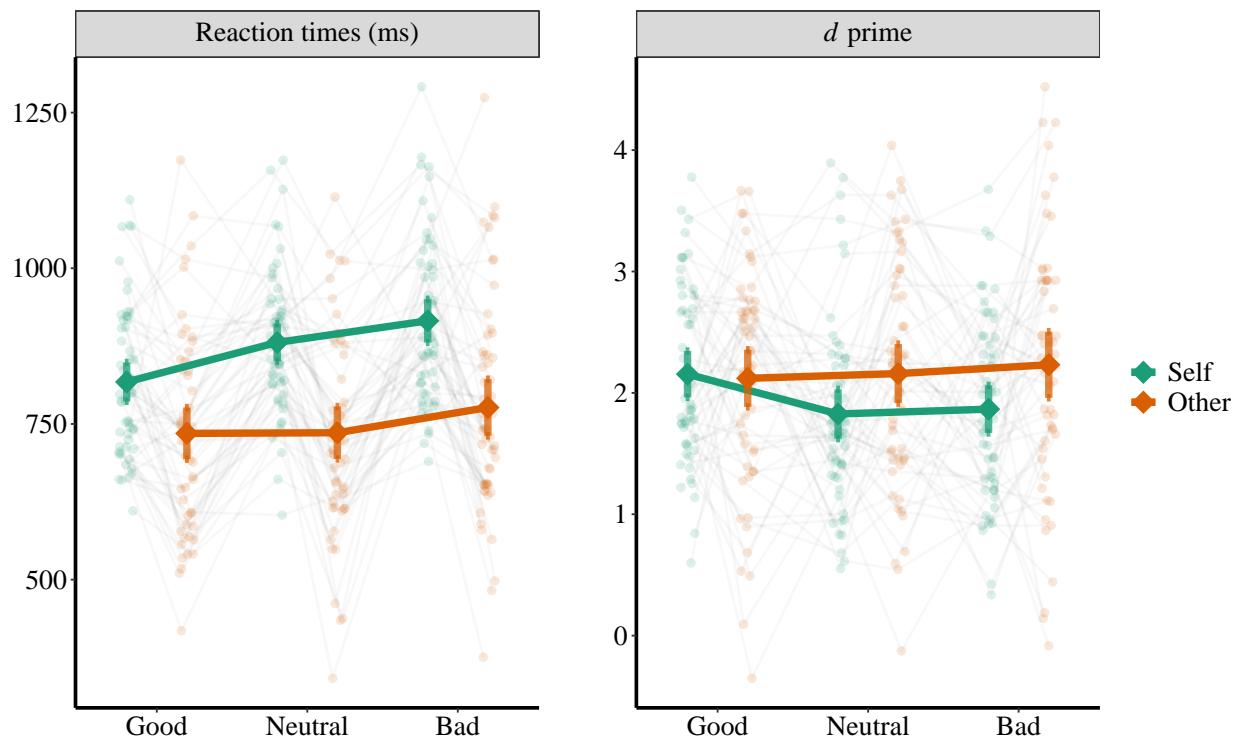


Figure 5. RT and *d* prime of Experiment 3b.

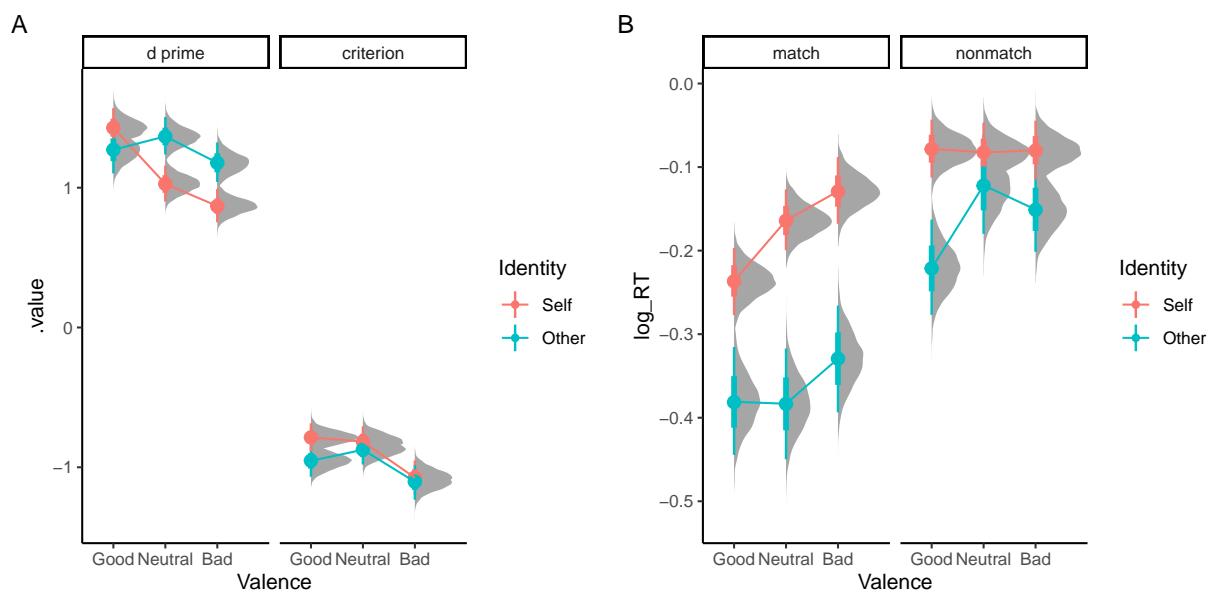


Figure 6. exp3b: Results of Bayesian GLM analysis.

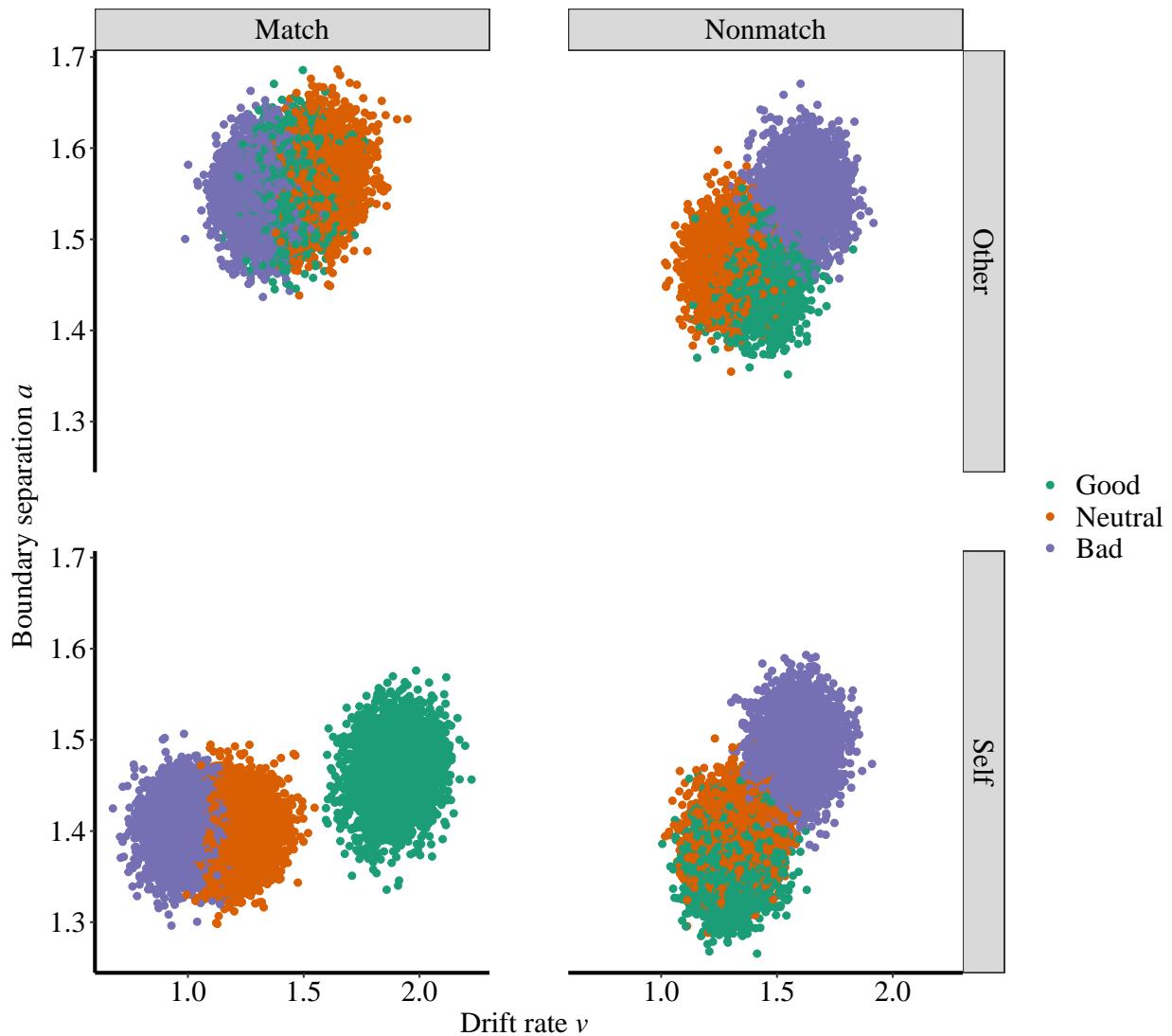


Figure 7. exp3b: Results of HDDM.

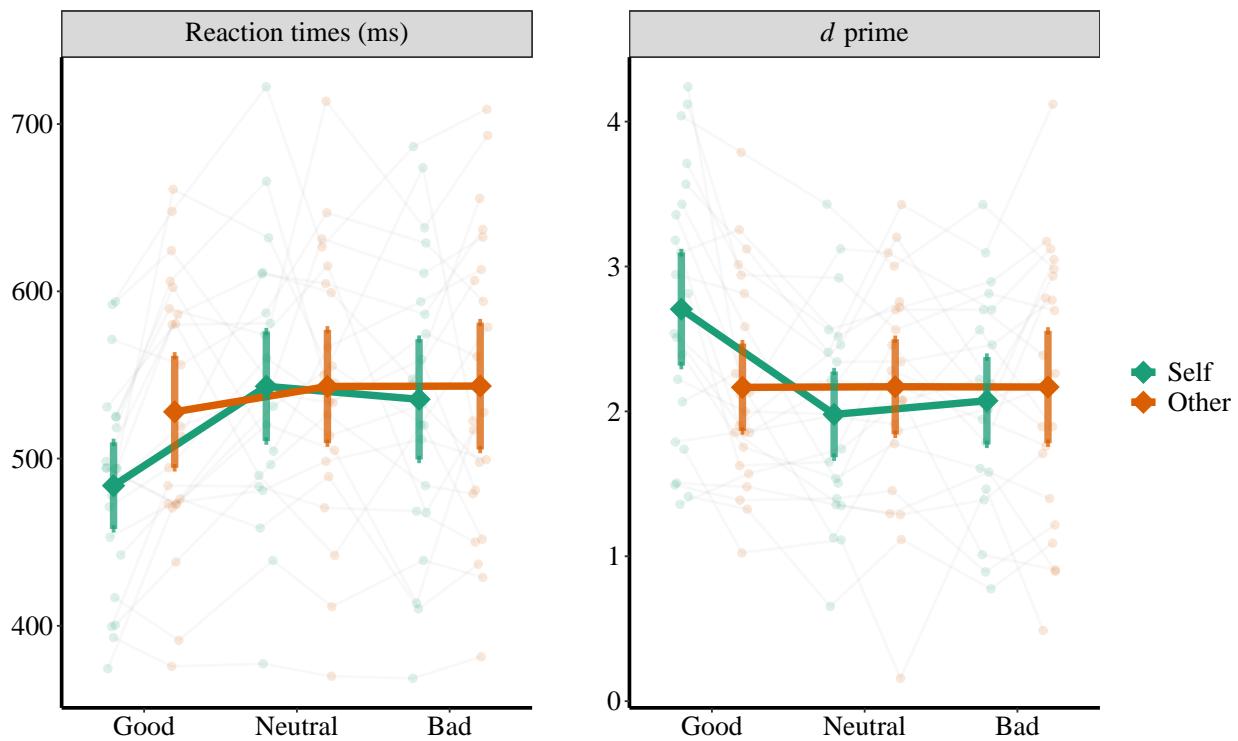


Figure 8. RT and  $d'$  of Experiment 6b.

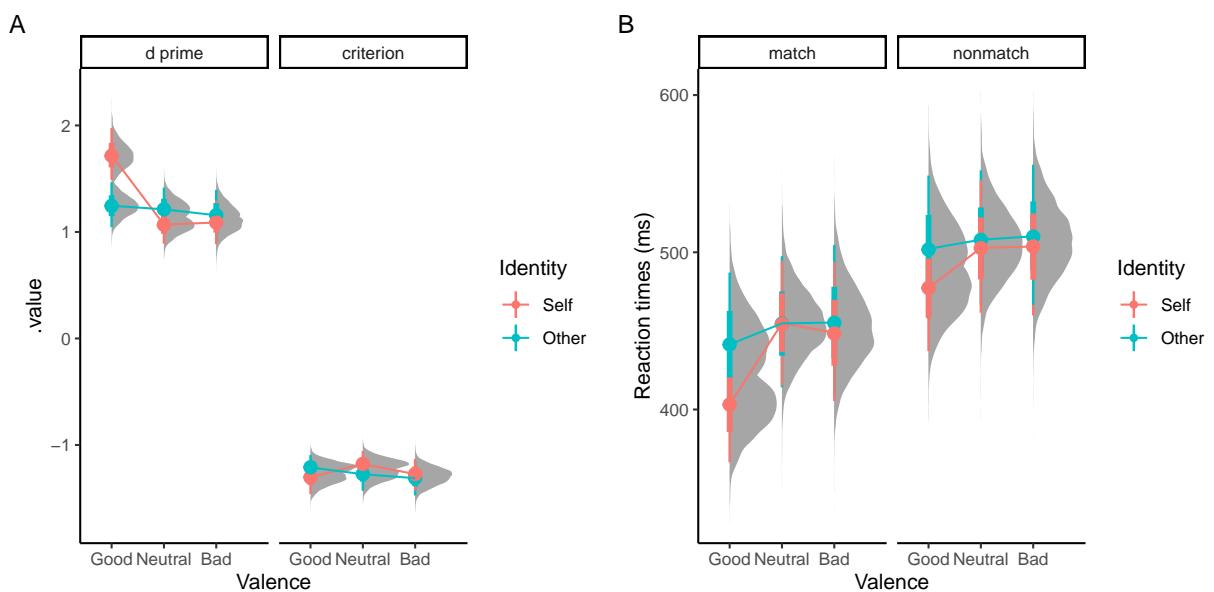


Figure 9. exp6b\_d1: Results of Bayesian GLM analysis.

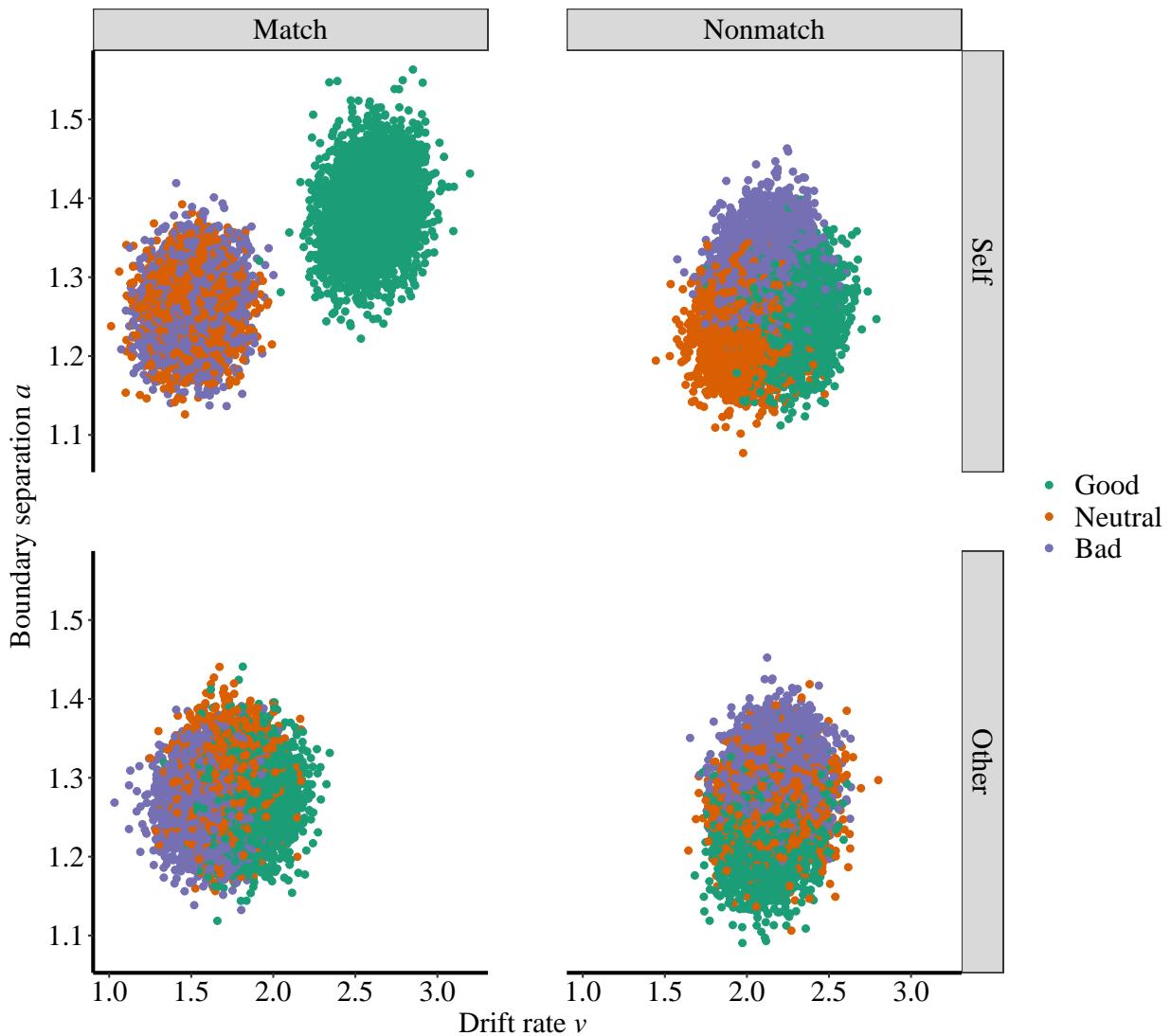


Figure 10. exp6b: Results of HDDM (Day 1).

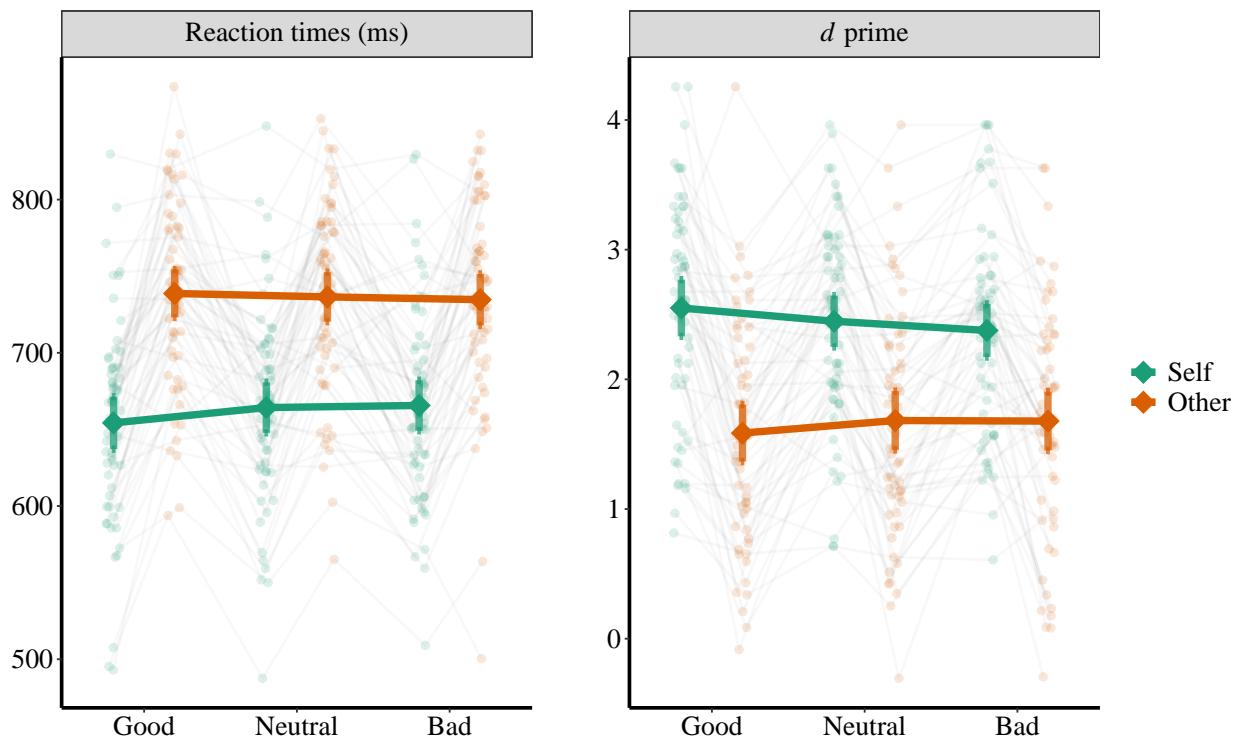


Figure 11. RT and  $d'$  of Experiment 4a.

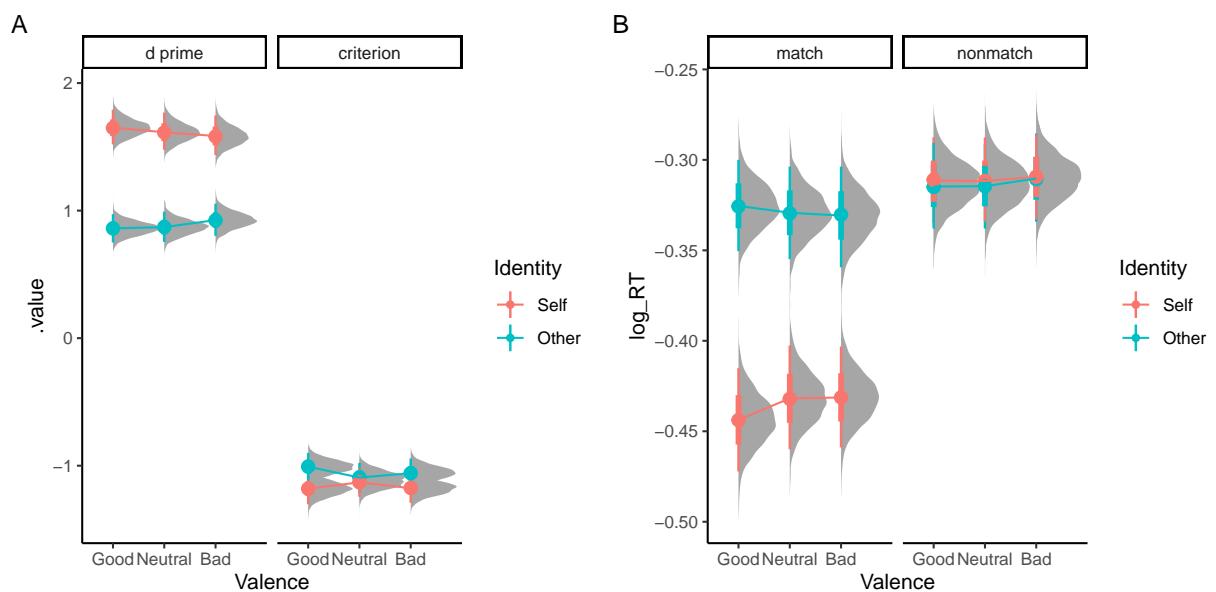


Figure 12. exp4a: Results of Bayesian GLM analysis.

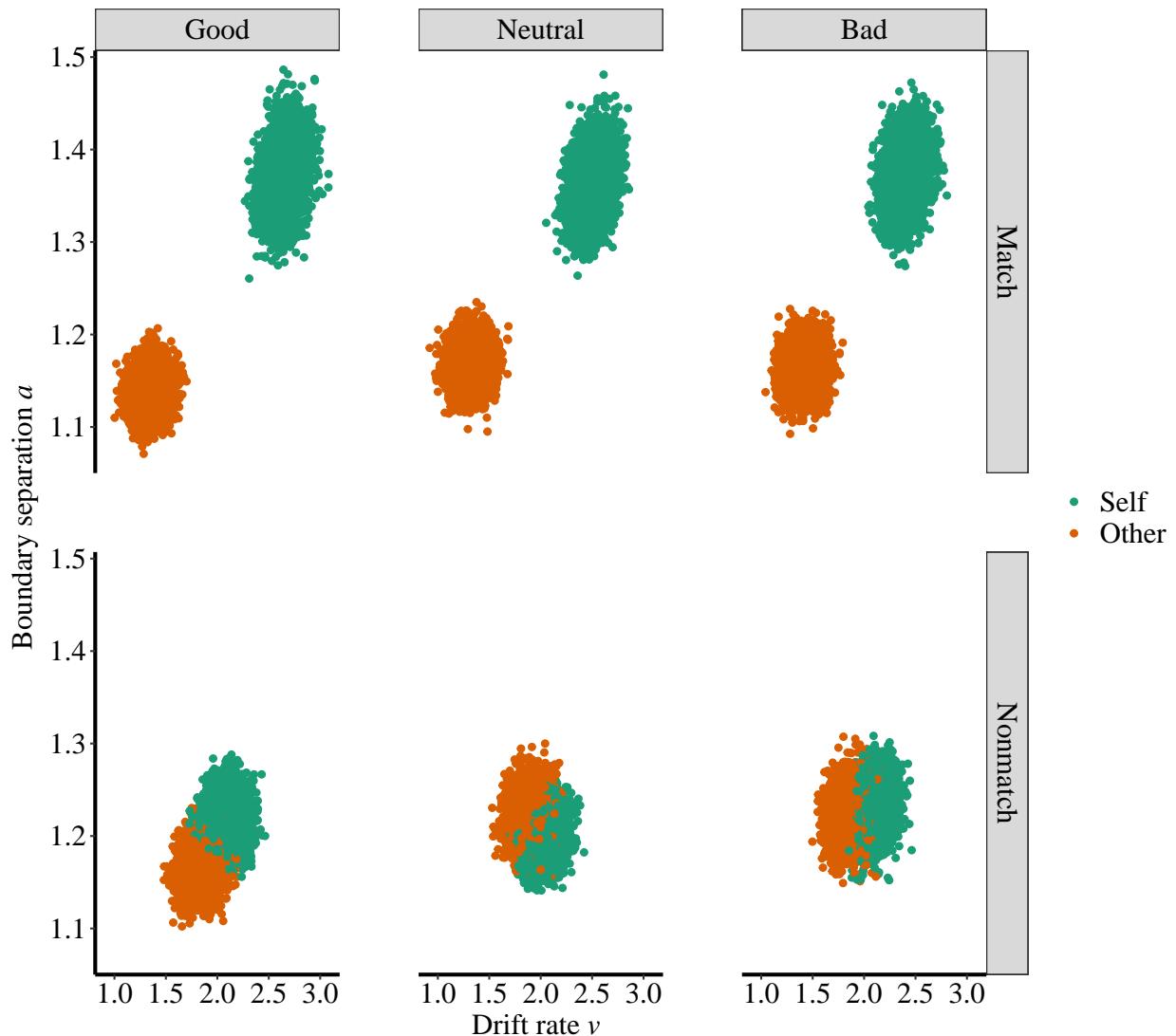


Figure 13. exp4a: Results of HDDM.

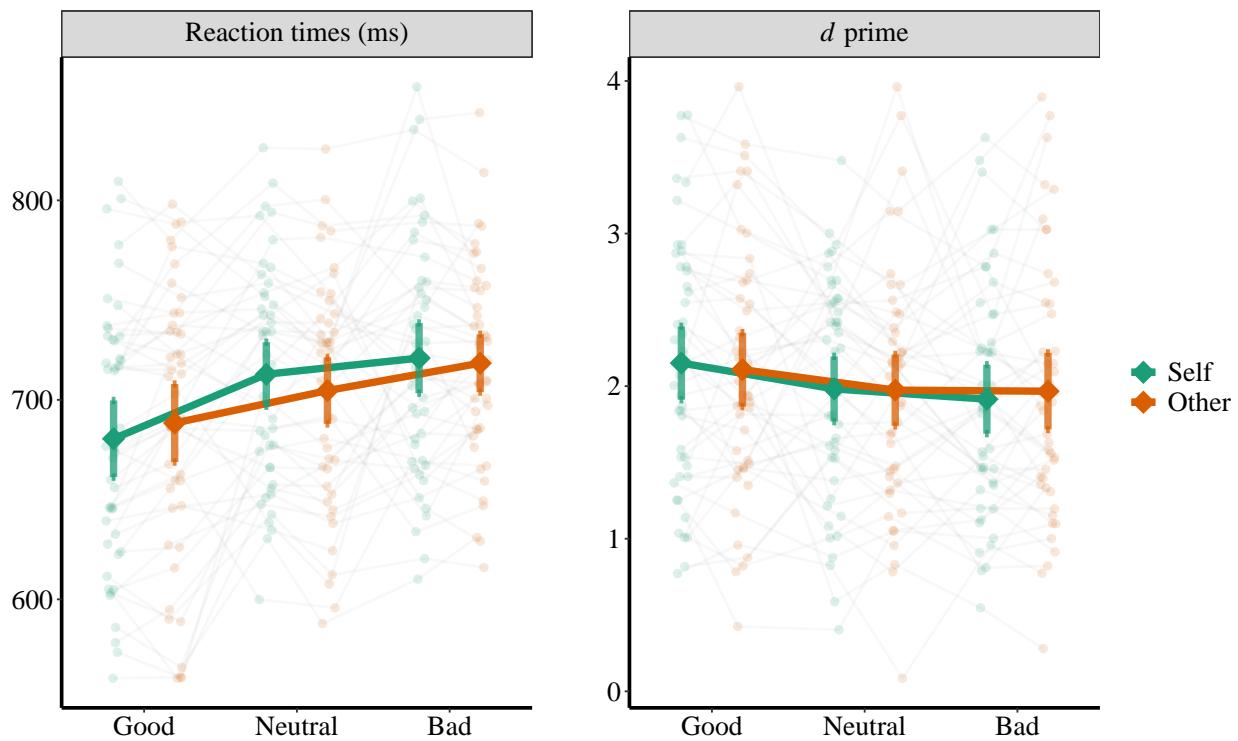


Figure 14. RT and  $d'$  prime of Experiment 4b.

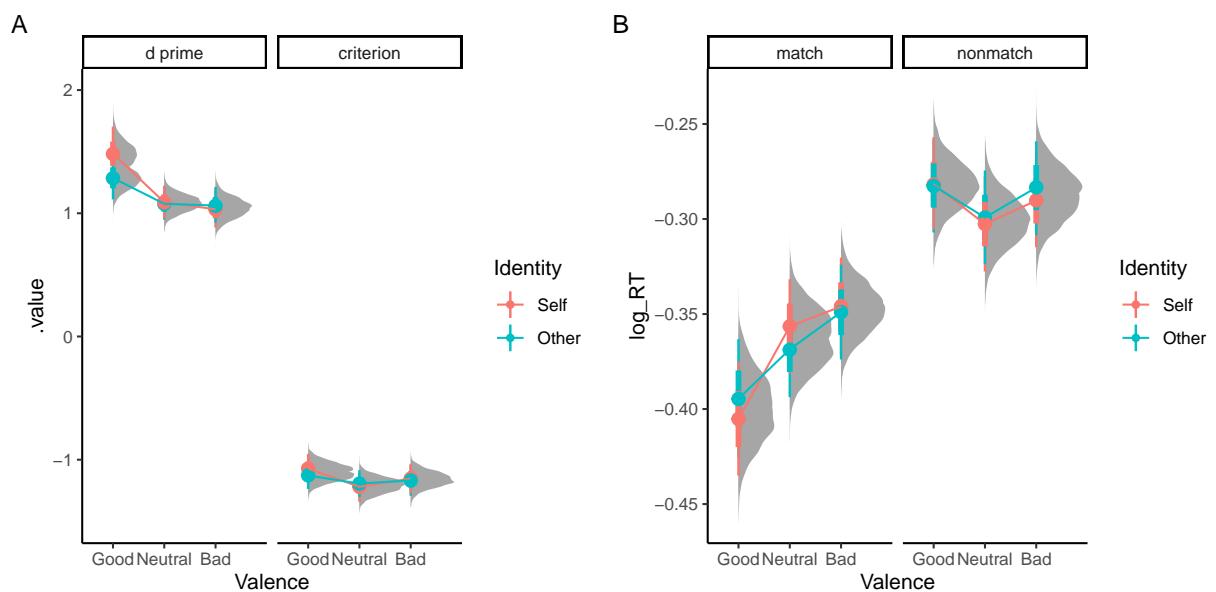


Figure 15. exp4b: Results of Bayesian GLM analysis.

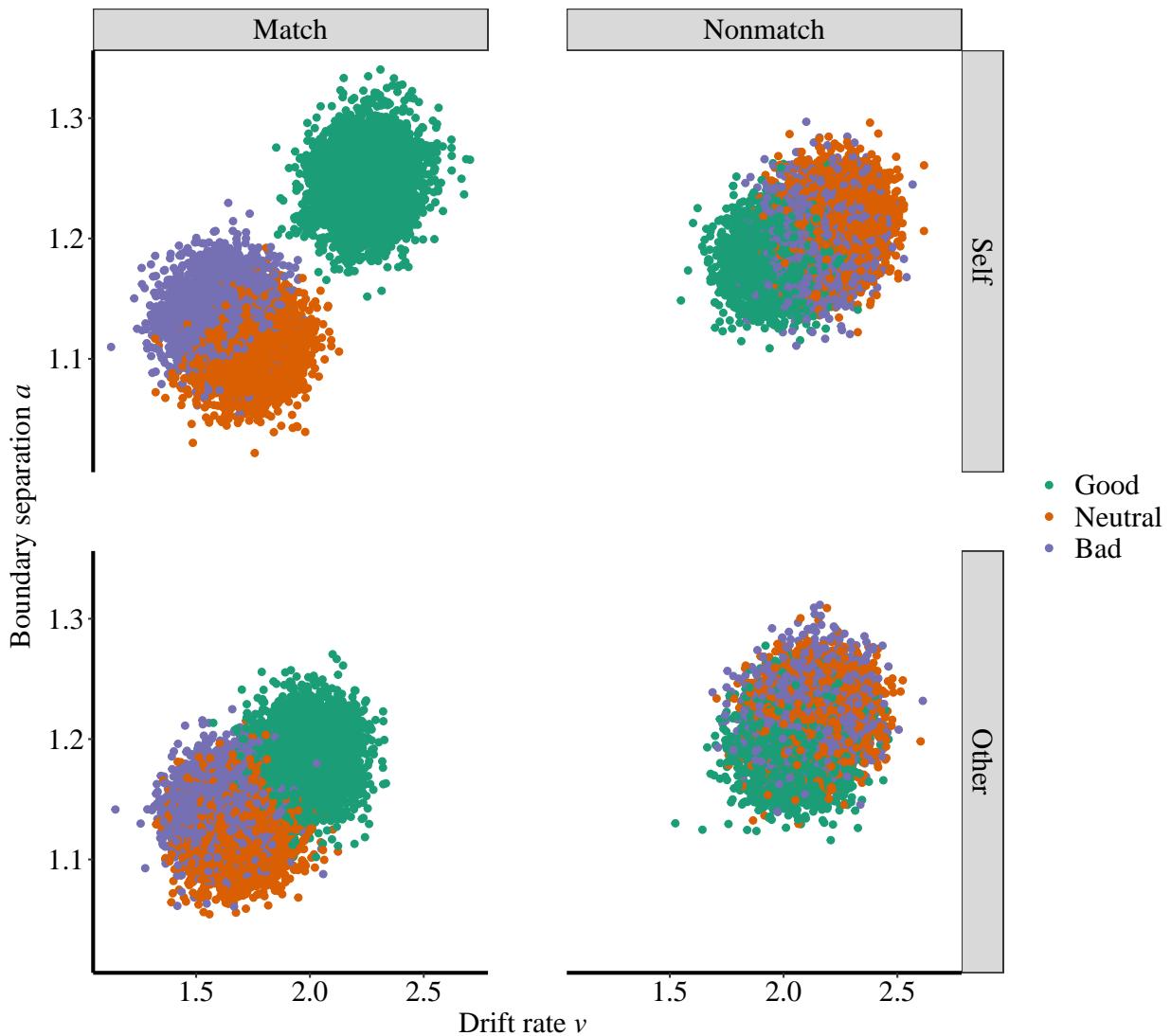


Figure 16. exp4b: Results of HDDM.

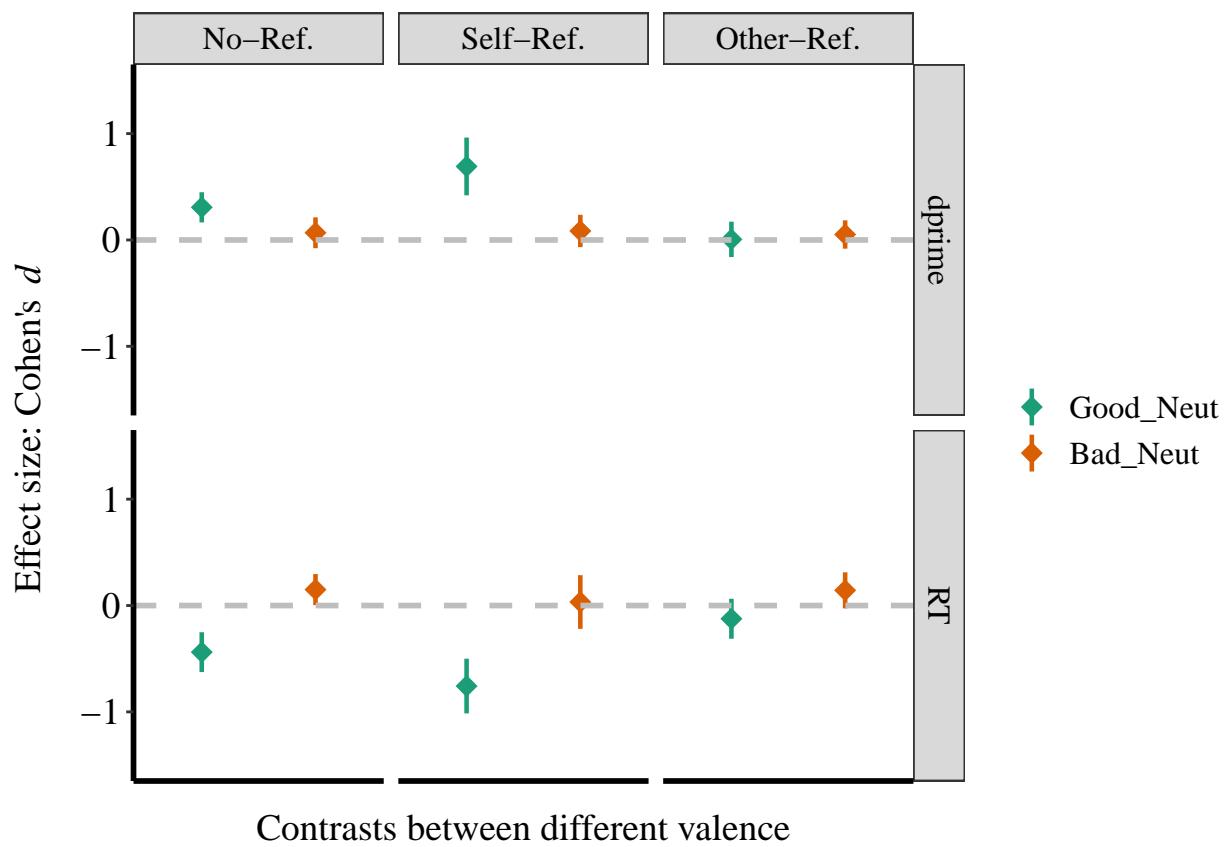


Figure 17. Effect size (Cohen's  $d$ ) of Valence.

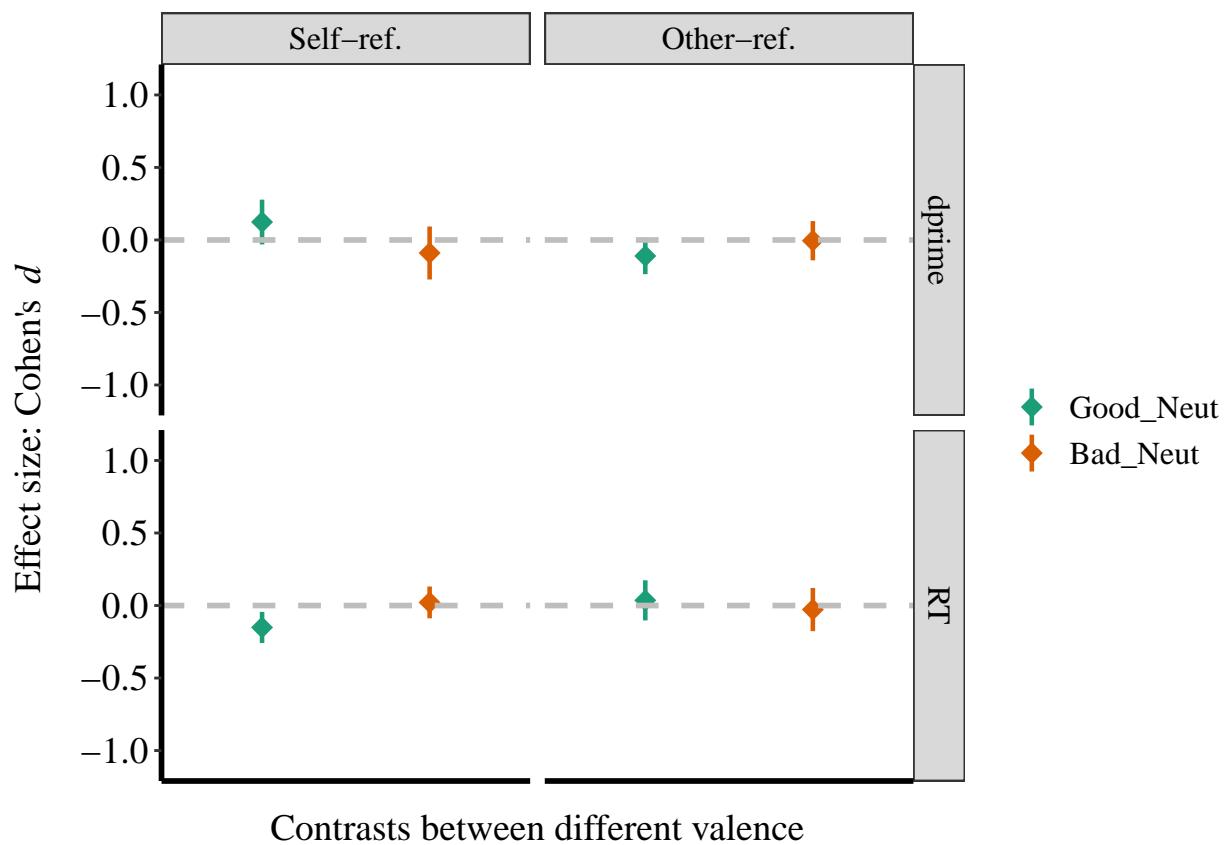


Figure 18. Effect size (Cohen's  $d$ ) of Valence in Exp4a.

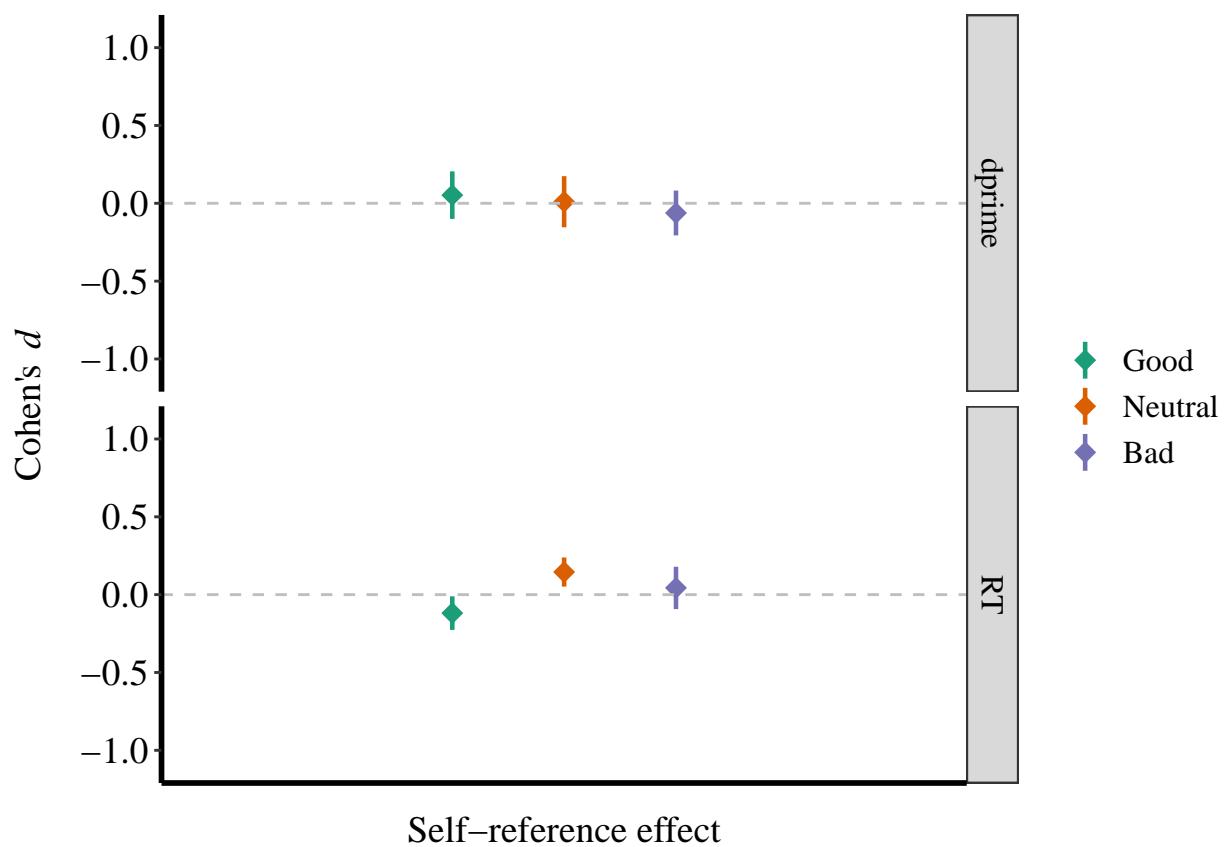


Figure 19. Effect size (Cohen's  $d$ ) of Valence in Exp4b.

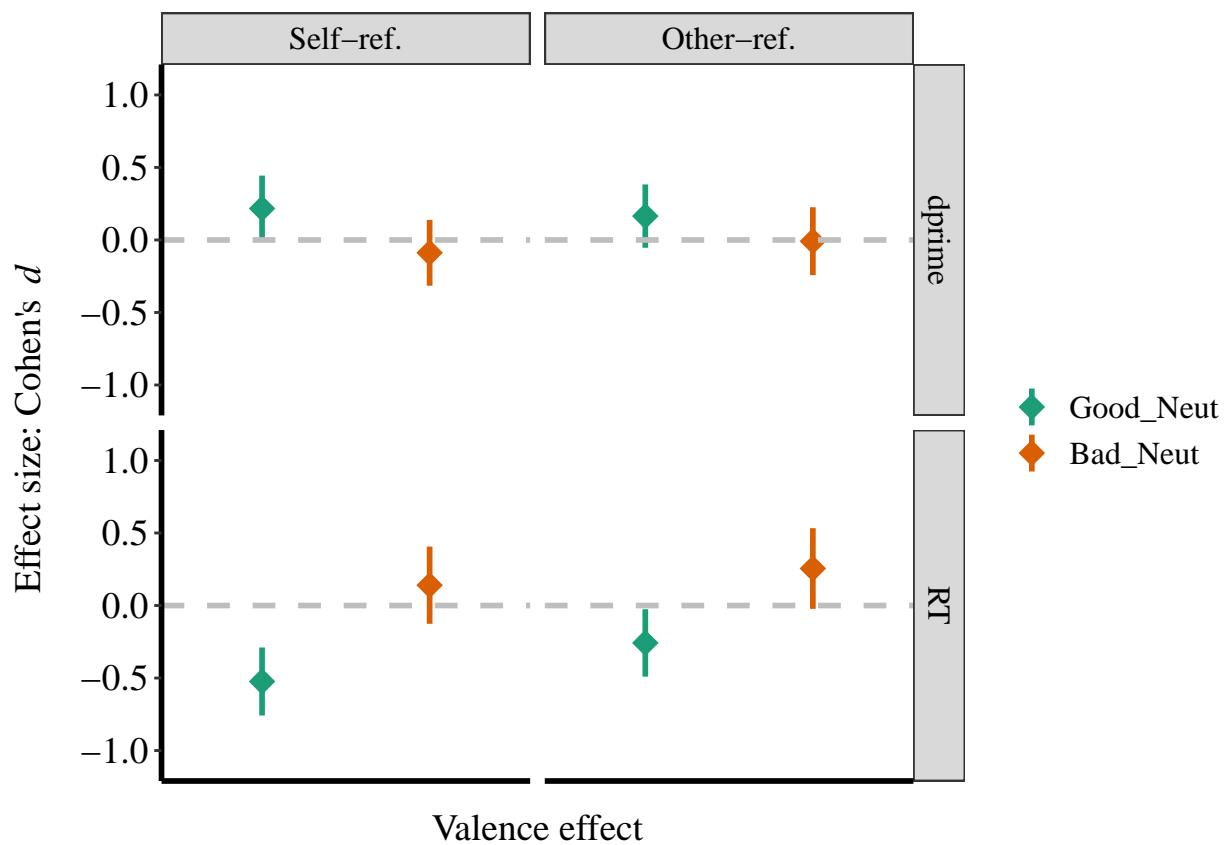


Figure 20. Effect size (Cohen's  $d$ ) of Valence in Exp4b.

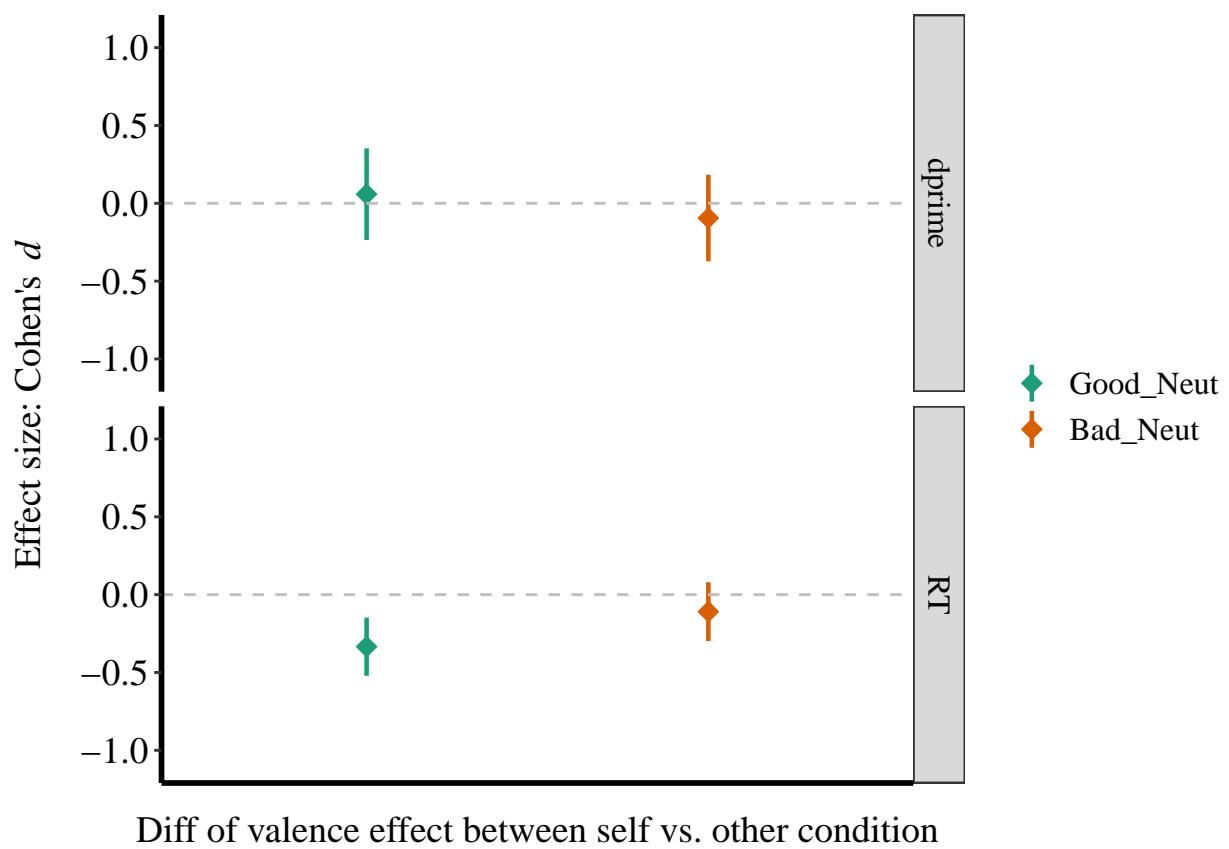


Figure 21. Effect size (Cohen's  $d$ ) of Valence in Exp4b.

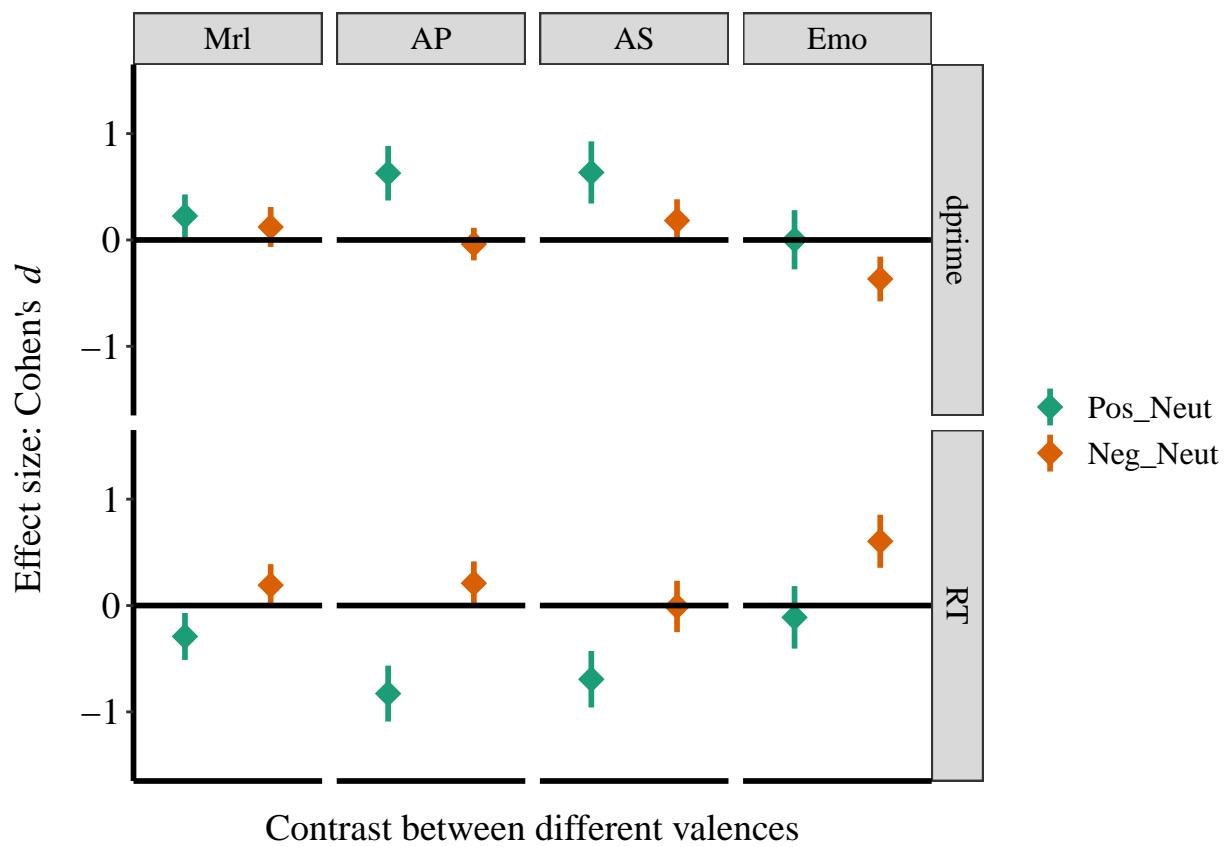


Figure 22. Effect size (Cohen's  $d$ ) of Valence in Exp5.

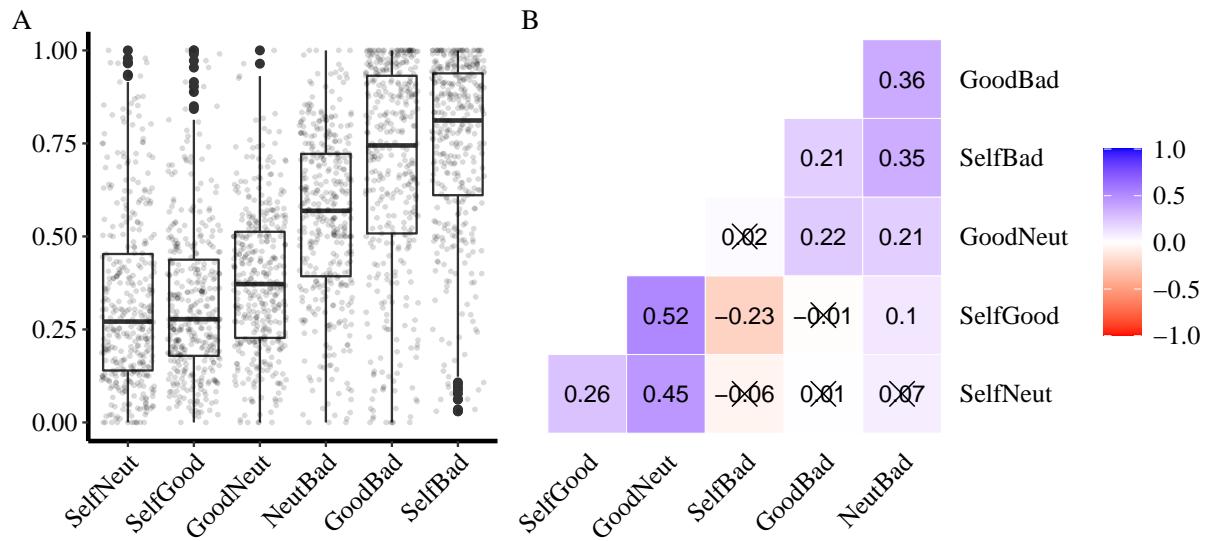


Figure 23. Self-rated personal distance

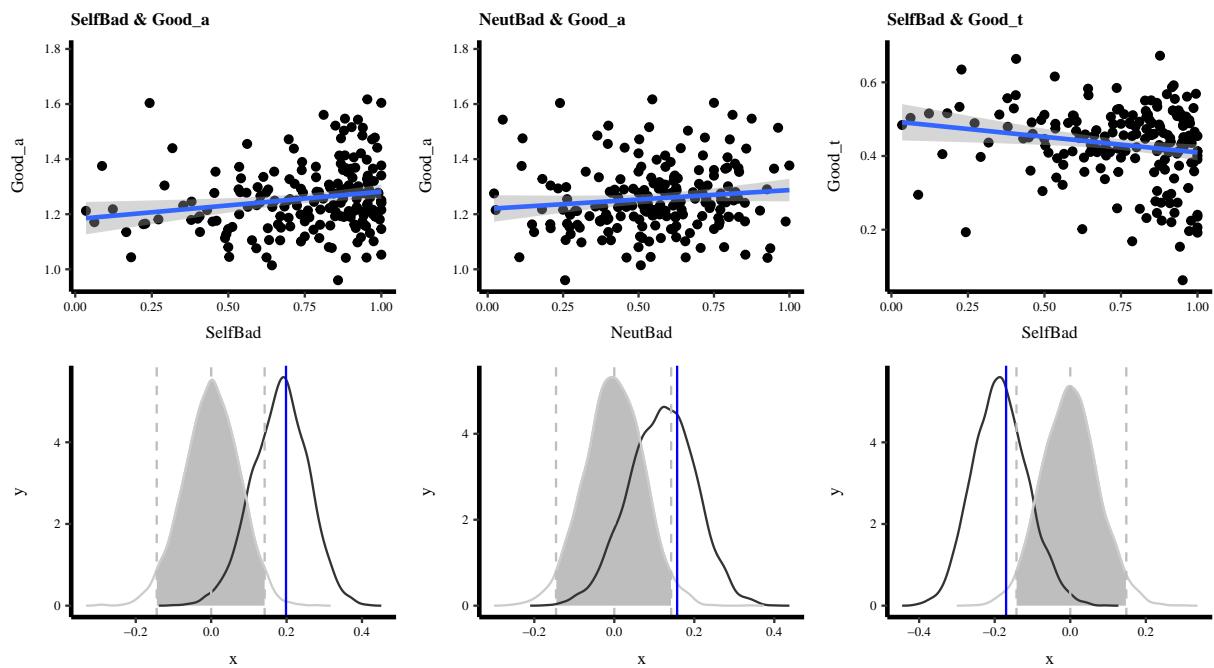
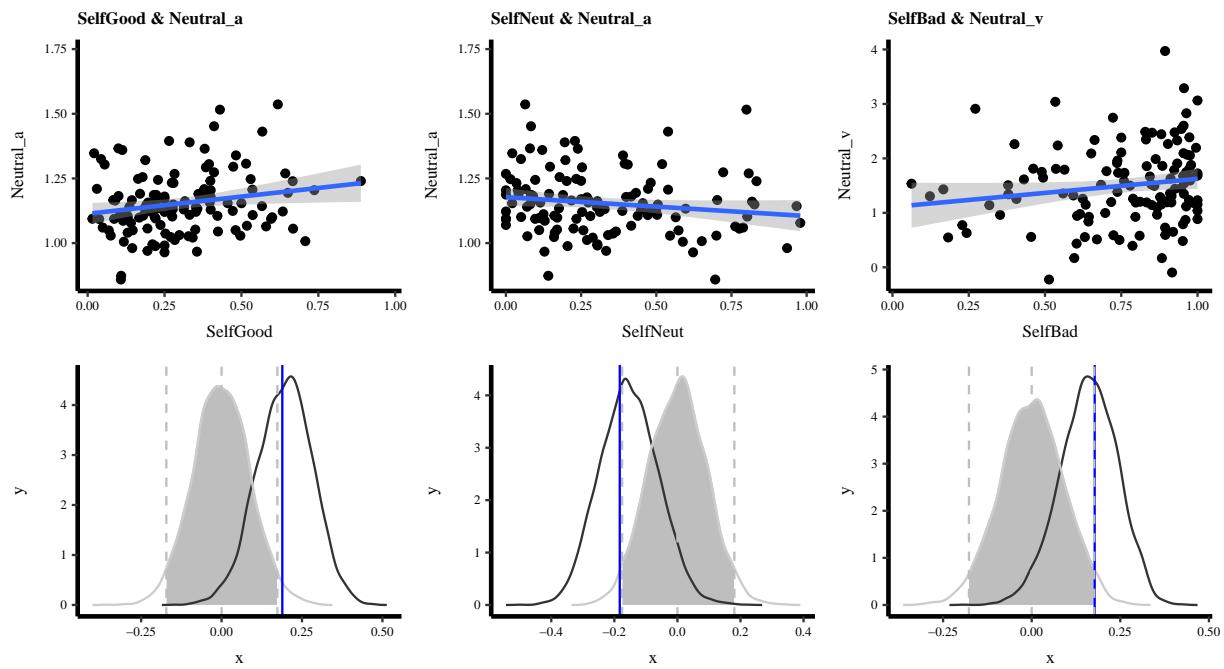


Figure 24. Correlation between moral identity and boundary separation of good condition; moral self-image and drift rate of good condition



*Figure 25.* Correlation between personal distance and boundary separation of neutral condition