

¹ Positivity bias in perceptual matching may reflect a spontaneous self-referential processing

² Hu Chuan-Peng^{1,2}, Kaiping Peng³, & Jie Sui^{3,4}

³ ¹ TBA

⁴ ² Leibniz Institute for Resilience Research, 55131 Mainz, Germany

⁵ ³ Tsinghua University, 100084 Beijing, China

⁶ ⁴ University of Aberdeen, Aberdeen, Scotland

⁷ Author Note

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

⁹ Department of Psychology, Tsinghua University, 100084 Beijing, China. Jie Sui, School of

¹⁰ Psychology, University of Aberdeen, Aberdeen, Scotland.

¹¹ Authors contriubtion: HCP, JS, & KP design the study, HCP collected the data,

¹² HCP analyzed the data and drafted the manuscript. KP & JS supported this project.

¹³ Correspondence concerning this article should be addressed to Hu Chuan-Peng,

¹⁴ Langenbeckstr. 1, Neuroimaging Center, University Medical Center Mainz, 55131 Mainz,

¹⁵ 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 make substantial amount of
⁴³ moral judgment in everyday life (e.g., Hofmann, Wisneski, Brandt, & Skitka, 2014). When
⁴⁴ making these judgments, people not only judge others behavior but, more importantly,
⁴⁵ others' character (Uhlmann, Pizarro, & Diermeier, 2015). That is, moral judgment is a way
⁴⁶ that we get information about other in the society, so that we can distinguish the good
⁴⁷ from the bad. This align with studies of person-perception in which moral character is a
⁴⁸ basic dimension of person perception in social life(Abele, Ellemers, Fiske, Koch, & Yzerbyt,
⁴⁹ 2020; Goodwin, 2015; Goodwin, Piazza, & Rozin, 2014; Willis & Todorov, 2006).

⁵⁰ **Morality as social categorization**

⁵¹ (McHugh, McGann, Igou, & Kinsella, 2019): Moral judgment as categorization

⁵² (MJAC); person-centered moral judgment person-perception

⁵³ People attribute the punishment to moral character instead of environmental factors

⁵⁴ (Dunlea & Heiphetz, 2020).

⁵⁵ **Self as an object of moral categorization: I am a good person**

⁵⁶ Social life: treat people the same way as we expect people to treat yourself. a person
⁵⁷ needs to both accurately evaluate others' moral character and behave in a way that she/he
⁵⁸ is perceived as a moral person, or at least not a morally bad person. Motivation to
⁵⁹ maintain good-self

60 **From good me to good us: Spontaneous good-self referential in moral
61 categorization**

62 Integrate self-categorization in moral judgment: an implicit good-me/good-we
63 assumption (Role of self in attractor model)

64 *self-categorization and variable self*

65 Given the importance of moral character, to successfully navigate in social world, a
66 person needs to both accurately evaluate others' moral character and behave in a way that
67 she/he is perceived as a moral person, or at least not a morally bad person. The former
68 was traditionally investigated as person perception in social psychology, while the latter
69 was studied separately as moral self concept (Monin & Jordan, 2009). As for the former
70 question, abundant of evidence revealed that people weigh morality heavily in evaluating
71 others (Goodwin, 2015) and evaluating the change of identity of others (Strohminger,
72 Knobe, & Newman, 2017). These findings suggest that morality has been internalized
73 standard for how people perceiving and remembering other people.

74 When it comes to self perception, there is accumulating evidence that people actively
75 maintain a good moral-self image. For example, recent research found that
76 self-enhancement effect is stronger than that in competence or sociability (Tappin &
77 McKay, 2017). Also, participants maintain their moral self-image even after their own
78 unethical behaviors (e.g., cheating)(Monin & Jordan, 2009). Similarly, when asked how
79 likely they will act ethically or unethically, most participants showed the tendency of less
80 likely to do unethical things (Klein & Epley, 2016). In other words, existing evidence
81 supported the notion that morality is important in person perception and self-concept,
82 people are motivated to maintain a good moral-self image.

83 [whether moral character information influence perception?, link to exp1a, b, c, and
84 exp2] However, as Freeman and Ambady (2011) put it, the focus of the studies in person
85 perception and moral self-concept is not to explain the perceptual process, rather, they are

86 explaining the higher-order social cognitive processes that come after. That is, current
87 studies on person perception are perception without perceptual process. Without
88 knowledge of perceptual processes, we can not have a full picture of how moral information
89 is processed in our cognition. As an increasing attention is paid to perceptual process
90 underlying social cognition, it's clear that perceptual processes are strongly influenced by
91 social factors, such as group-categorization, stereotype (see Xiao, Coppin, & Bavel, 2016;
92 Stolier & Freeman, 2016). Given the importance of moral character and that moral
93 character related information has strong influence on learning and memory (Carlson,
94 Maréchal, Oud, Fehr, & Crockett, 2020; Stanley & De Brigard, 2019), one might expect
95 that moral character related information could also play a role in perceptual process.

96 [using associative learning task to study the moral character's influence on perception]
97 Though theoretically possible, no empirical studies had directly addressed this issue. There
98 were only a few studies about the temporal dynamics of judging the trustworthiness of face
99 (e.g., Dzhelyova, Perrett, & Jentzsch, 2012), but trustworthy is not equal to morality. One
100 difficulty of studying moral character's influence on perceptual decision-making is that
101 moral character is a high-level and hidden state instead of observable feature. usually, one
102 needs necessary information, e.g., behavior history, to infer moral character of a person.
103 For example, Anderson et al. (2011) asked participant to first study the behavioral
104 description of faces and then asked participant to perform a perceptual detection task.
105 They assumed that by learning the behavioral description of a person (represented by a
106 face), participants could acquire the moral related information about that face, and this
107 association would then bias the perceptual processing of the faces (but see Stein, Grubb,
108 Bertrand, Suh, and Verosky (2017)). However, one drawback of this approach is that
109 participants may differ greatly when inferring the moral character of the person from
110 behavioral descriptions, given that notion what is morality itself is varying across
111 population (Henrich, Heine, & Norenzayan, 2010) and those descriptions and faces may
112 themselves are idiosyncratic, therefore, introduced large variation in experimental design.

An alternative is to use abstract semantic concepts to study how these concepts influence perception. Abstract concepts of moral character is part of our daily life and it can be used to represent rich information. These abstract concepts is part of a dynamic network in which sensory cue, concrete behaviors and other information can activate/inhibit each other (e.g., aggressiveness) (Amodio, 2019; Freeman & Ambady, 2011). If a concept of moral character (e.g., good person) is activated, it should also be able to influence on the perceptual process of the visual cues through the dynamic network, especially when the perceptual decision-making is about the concept-cue association. In this case, abstract concepts of moral character may serve as signal of moral reputation (for others) or moral self-concept. Indeed, previous studies used the moral words and found that moral related information can be perceived faster (Gantman & Van Bavel, 2014, but see @firestone_enhanced_2015). If moral character is an important in person perception, then, just as those other information such as races, education, (see Xiao et al., 2016), moral character related information might change the perceptual processes.

To investigate the above possibility, we used an associative learning paradigm to study how moral character concept change perceptual decision-making. In this paradigm, simple geometric shapes were paired with different words whose dominant use is to describe the moral character of person. Participants first learn the associations between shapes and words, i.e. triangle is a good-person, i.e., building direct association between high-level, hidden moral characters and visual cues. After remembered the associations, they perform a matching task to judge whether the shape-word pair presented on the screen match the association they learned. This paradigm has been used in studying the perceptual process of self-concept, but had also proven useful in studying other concepts like social group (Enock, Sui, Hewstone, & Humphreys, 2018). By using simple and morally neutral shapes, we controlled the variations caused by visual cues.

Our first question is, whether the words used the in the associative paradigm is really related to the moral character? As we reviewed above, previous theories, especially the

140 interactive dynamic theory, would support this assumption. To validate that moral
141 character concepts activated moral character as a social cue, we used four experiments to
142 explore and validate the paradigm. The first experiment direct adopted associative
143 paradigm and change the words from “self”, “other” to “good-person”, “neutral-person”,
144 and “bad-person”. Then, we change the words to the ones that have more explicit moral
145 meaning (“kind-person”, “neutral-person”, and “evil-person”). Then, as in Anderson et al.
146 (2011), we asked participant to learn the behavioral history of three different names, and
147 then use the names, as moral character words, for associative learning. Finally, we also
148 tested that simultaneously present shape-word pair and sequentially present word and
149 shape didn’t change the pattern. All of these four experiments showed a robust effect of
150 moral character, that is, the positive moral character associated stimuli were prioritized.

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

¹⁶⁷ McGarty, 1994). However, the above four experiments can not distinguish between these
¹⁶⁸ two possibilities, because there are evidence for both reward- (e.g., Sui, He, & Humphreys,
¹⁶⁹ 2012) and in-group (Enock et al., 2018) prioritization.

¹⁷⁰ [Distinguish two explanations by make self salient, exp3a, 3b, 6b] Though these two
¹⁷¹ frameworks can both account for the positivity effect found in first four experiments (i.e.,
¹⁷² prioritization of “good-person”, but not “neutral person” and “bad person”), they have
¹⁷³ different prediction if the experiment design include both identity and moral valence where
¹⁷⁴ the valence (good, bad, and neutral) conditions can describe both self and other. In this
¹⁷⁵ case the identity become salient and participants are less likely to spontaneously identify a
¹⁷⁶ good-person other than self as in-group, but the value of good-person still exists. This
¹⁷⁷ means that the social categorization theory predicts participants prioritize good-self but not
¹⁷⁸ good-other, while reward-based attention theory predicts participants are both prioritized.
¹⁷⁹ Also, as in Hu, Lan, Macrae, and Sui (2020), people may also only identify with good-self
¹⁸⁰ instead of bad self. That is, people will show a unique pattern of self-identification: only
¹⁸¹ good-self is identified as “self” while all the others categories were excluded.

¹⁸² In exp 3a, 3b, and 6b, we found that (1) good-self is always faster than neutral-self
¹⁸³ and bad-self, but good-other only have weak to null advantage to neutral-other and
¹⁸⁴ bad-other. which mean the social categorization is self-centered. (2) good-self’s advantage
¹⁸⁵ over good other only occur when self- and other- were in the same task. i.e. the relative
¹⁸⁶ advantage is competition based instead of absolute. These three experiments suggest that
¹⁸⁷ people more like to view the moral character stimuli as person and categorize good-self as
¹⁸⁸ an unique category against all others. A mini-meta-analysis showed that there was no
¹⁸⁹ effect of valence when the identity is other.

¹⁹⁰ [what we care? valence of the self exp4a or identity of the good exp4b?] Next, we go
¹⁹¹ further to disentangle the good-self complex: people care more about whether the self is
¹⁹² good, or whether the good is self, or both? If people care more about whether the self is

193 good, then, subtle cue of the valence may have an impact on perceptual process of the self.
194 In contrast, if people care more about whether the good's identity, i.e., whether the good is
195 self, then subtle cue of identity (self v. other) may have a impact on perceptual process of
196 the good. We tested the good-self complex with two more experiments. In exp 4a (id is
197 task-relevant, valence is task-irrelevant), if people care about the valence of the self, then,
198 the task-irrelevant information may influence the processing of the self. While in exp 4b
199 (valence is task-relevant, id is task-irrelevant), if people care about the id of the good, then,
200 the task-irrelevant id information will has a influence on the process of the good.

201 [Good self in self-reported data] As an exploration, we also collected participants'
202 self-reported psychological distance between self and good-person, bad-person, and
203 neutral-person. We explored the general pattern and the correlation between self-reported
204 distance and reaction-based indices.

205 [whether categorize self as positive is not limited to morality] Finally, we explored the
206 pattern is generalized to all positive traits or only to morality. We found that
207 self-categorization is not limited to morality, even morality is central to social life. we used
208 aesthetic aspect as another instance.

209 Key concepts and discussing points: *Self-categories* are cognitive groupings of self and
210 some class of stimuli as identical or different from some other class. [Turner et al.]
211 *Personal identity* refers to self-categories that define the individual as a unique person in
212 terms of his or her individual differences from other (in-group) persons. *Social identity*
213 refers to the shared social categorical self ("us" vs. "them"). *variable self*: Who we are,
214 how we see ourselves, how we define our relations to others (indeed whether they are
215 construed as "other" or as part of the extended "we" self) is different in different settings.
216 *Identification*: the degree to which an individual feels connected to an ingroup or includes
217 the ingroup in his or her self-concept. (self is not bad;) Morality as a way for
218 social-categorization? People are more likely to identify themselves with trustworthy faces

²¹⁹ (Verosky & Todorov, 2010) (trustworthy faces has longer RTs). What is the relation
²²⁰ between morally good and self in a semantic network (attractor network) (Freeman &
²²¹ Ambady, 2011). How to deal with the *variable self* (self-categorization theory)
²²² vs. *core/true/authentic self* vs. *self-enhancement*

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

Disclosures

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

230 All the experiments reported were not pre-registered. Most experiments (1a ~ 6b,
231 except experiment 3b) reported in the current study were first finished between 2014 to
232 2016 in Tsinghua University, Beijing, China. Participants in these experiments were
233 recruited in the local community. To increase the sample size of experiments to 50 or more
234 (Simmons, Nelson, & Simonsohn, 2013), we recruited additional participants in Wenzhou
235 University, Wenzhou, China in 2017 for experiment 1a, 1b, 4a, and 4b. Experiment 3b was
236 finished in Wenzhou University in 2017. To have a better estimation of the effect size, we
237 included the data from two experiments (experiment 7a, 7b) that were reported in Hu et
238 al. (2020) (See Table S1 for overview of these experiments).

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

241

General methods

242 **Design and Procedure**

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

255 The experiment 1a, 1b, 1c, 2, and 6a shared a 2 (matching: match vs. nonmatch) by
256 3 (moral valence: good vs. neutral vs. bad) within-subject design. Experiment 1a was the
257 first one of the whole series studies and 1b, 1c, and 2 were conducted to exclude the
258 potential confounding factors. More specifically, experiment 1b used different Chinese
259 words as label to test whether the effect only occurred with certain familiar words.
260 Experiment 1c manipulated the moral valence indirectly: participants first learned to
261 associate different moral behaviors with different neutral names, after remembered the
262 association, they then performed the perceptual matching task by associating names with
263 different shapes. Experiment 2 further tested whether the way we presented the stimuli
264 influence the effect of valence, by sequentially presenting labels and shapes. Note that part
265 of participants of experiment 2 were from experiment 1a because we originally planned a
266 cross task comparison. Experiment 6a, which shared the same design as experiment 2, was

267 an EEG experiment which aimed at exploring the neural correlates of the effect. But we
268 will focus on the behavioral results of experiment 6a in the current manuscript.

269 For experiment 3a, 3b, 4a, 4b, 6b, 7a, and 7b, we included self-reference as another
270 within-subject variable in the experimental design. For example, the experiment 3a directly
271 extend the design of experiment 1a into a 2 (matchness: match vs. nonmatch) by 2
272 (reference: self vs. other) by 3 (moral valence: good vs. neutral vs. bad) within-subject
273 design. Thus in experiment 3a, there were six conditions (good-self, neutral-self, bad-self,
274 good-other, neutral-other, and bad-other) and six shapes (triangle, square, circle, diamond,
275 pentagon, and trapezoids). The experiment 6b was an EEG experiment extended from
276 experiment 3a but presented the label and shape sequentially. Because of the relatively
277 high working memory load (six label-shape pairs), experiment 6b were conducted in two
278 days: the first day participants finished perceptual matching task as a practice, and the
279 second day, they finished the task again while the EEG signals were recorded. Experiment
280 3b was designed to separate the self-referential trials and other-referential trials. That is,
281 participants finished two different blocks: in the self-referential blocks, they only responded
282 to good-self, neutral-self, and bad-self, with half match trials and half non-match trials; for
283 the other-reference blocks, they only responded to good-other, neutral-other, and
284 bad-other. Experiment 7a and 7b were designed to test the cross task robustness of the
285 effect we observed in the aforementioned experiments (see, Hu et al., 2020). The matching
286 task in these two experiments shared the same design with experiment 3a, but only with
287 two moral valence, i.e., good vs. bad. We didn't include the neutral condition in
288 experiment 7a and 7b because we found that the neutral and bad conditions constantly
289 showed non-significant results in experiment 1 ~ 6.

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

294 4b, we reversed the role of self-reference and valence in the task: participant learnt three
295 labels (good-person, neutral-person, and bad-person) and three shapes (circle, square, and
296 triangle), and the words related to identity, “self” or “other”, were presented in the shapes.
297 As in 4a, participants were told to ignore the words inside the shape during the task.

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

304 E-prime 2.0 was used for presenting stimuli and collecting behavioral responses,
305 except that experiment 7a and 7b used Matlab Psychtoolbox (Brainard, 1997; Pelli, 1997).
306 For participants recruited in Tsinghua University, they finished the experiment individually
307 in a dim-lighted chamber, stimuli were presented on 22-inch CRT monitors and their head
308 were fixed by a chin-rest brace. The distance between participants’ eyes and the screen was
309 about 60 cm. The visual angle of geometric shapes was about $3.7^\circ \times 3.7^\circ$, the fixation cross
310 is of ($0.8^\circ \times 0.8^\circ$ of visual angle) at the center of the screen. The words were of $3.6^\circ \times 1.6^\circ$
311 visual angle. The distance between the center of the shape or the word and the fixation
312 cross was 3.5° of visual angle. For participants recruited in Wenzhou University, they
313 finished the experiment in a group consisted of 3 ~ 12 participants in a dim-lighted testing
314 room. Participants were required to finished the whole experiment independently. Also,
315 they were instructed to start the experiment at the same time, so that the distraction
316 between participants were minimized. The stimuli were presented on 19-inch CRT monitor.
317 The visual angles are could not be exactly controlled because participants’s chin were not
318 fixed.

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

³²⁰ questionnaire after they finish the behavioral tasks. All the questionnaire data are open
³²¹ (see, dataset 4 in Liu et al., 2020). See Table S1 for a summary information about all the
³²² experiments.

³²³ **Data analysis**

³²⁴ **Analysis of individual study.** We used the `tidyverse` of r (see script
³²⁵ `Load_save_data.r`) to exclude the practicing trials, invalid trials of each participants, and
³²⁶ invalid participants, if there were any, in the raw data. Results of each experiment were
³²⁷ then analyzed in three different approaches.

³²⁸ ***Classic NHST.***

³²⁹ First, as in Sui et al. (2012), we analyzed the accuracy and reaction times using
³³⁰ classic repeated measures ANOVA in the Null Hypothesis Significance Test (NHST)
³³¹ framework. Repeated measures ANOVAs is essentially a two-step mixed model. In the first
³³² step, we estimate the parameter on individual level, and in the second step, we used
³³³ repeated ANOVA to test the Null hypothesis. More specifically, for the accuracy, we used a
³³⁴ signal detection approach, in which individual' sensitivity d' was estimated first. To
³³⁵ estimate the sensitivity, we treated the match condition as the signal while the nonmatch
³³⁶ conditions as noise. Trials without response were coded either as “miss” (match trials) or
³³⁷ “false alarm” (nonmatch trials). Given that the match and nonmatch trials are presented
³³⁸ in the same way and had same number of trials across all studies, we assume that
³³⁹ participants' inner distribution of these two types of trials had equal variance but may had
³⁴⁰ different means. That is, we used the equal variance Gaussian SDT model (EVSDT) here
³⁴¹ (Rouder & Lu, 2005). The d' was then estimated as the difference of the standardized hit
³⁴² and false alarm rats (Stanislaw & Todorov, 1999):

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

³⁴³ where the HR means hit rate and the FAR mean false alarm rate. zHR and $zFAR$ are

³⁴⁴ the standardized hit rate and false alarm rates, respectively. These two z -scores were
³⁴⁵ converted from proportion (i.e., hit rate or false alarm rate) by inverse cumulative normal
³⁴⁶ density function, Φ^{-1} (Φ is the cumulative normal density function, and is used convert z
³⁴⁷ score into probabilities). Another parameter of signal detection theory, response criterion c ,
³⁴⁸ is defined by the negative standardized false alarm rate (DeCarlo, 1998): $-zFAR$.

³⁴⁹ For the reaction times (RTs), only RTs of accurate trials were analyzed. We first
³⁵⁰ calculate the mean RTs of each participant and then subject the mean RTs of each
³⁵¹ participant to repeated measures ANOVA. Note that we set the alpha as .05. The repeated
³⁵² measure ANOVA was done by `afex` package (<https://github.com/singmann/afex>).

³⁵³ To control the false positive rate when conducting the post-hoc comparisons, we used
³⁵⁴ Bonferroni correction.

³⁵⁵ ***Bayesian hierarchical generalized linear model (GLM).***

³⁵⁶ The classic NHST approach may ignore the uncertainty in estimate of the parameters
³⁵⁷ for SDT (Rouder & Lu, 2005), and using mean RT assumes normal distribution of RT
³⁵⁸ data, which is always not true because RTs distribution is skewed (Rousselet & Wilcox,
³⁵⁹ 2019). To better estimate the uncertainty and use a more appropriate model, we also tried
³⁶⁰ Bayesian hierarchical generalized linear model to analyze each experiment's accuracy and
³⁶¹ RTs data. We used BRMs (Bürkner, 2017) to build the model, which used Stan (Carpenter
³⁶² et al., 2017) to estimate the posterior.

³⁶³ In the GLM model, we assume that the accuracy of each trial is Bernoulli distributed
³⁶⁴ (binomial with 1 trial), with probability p_i that $y_i = 1$.

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

³⁶⁵ In the perceptual matching task, the probability p_i can then be modeled as a function of
³⁶⁶ the trial type:

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

367 The outcomes y_i are 0 if the participant responded “nonmatch” on trial i , 1 if they
 368 responded “match”. The probability of the “match” response for trial i for a participant is
 369 p_i . We then write the generalized linear model on the probits (z-scores; Φ , “Phi”) of ps . Φ
 370 is the cumulative normal density function and maps z scores to probabilities. Given this
 371 parameterization, the intercept of the model (β_0) is the standardized false alarm rate
 372 (probability of saying 1 when predictor is 0), which we take as our criterion c . The slope of
 373 the model (β_1) is the increase of saying 1 when predictor is 1, in z -scores, which is another
 374 expression of d' . Therefore, $c = -zHR = -\beta_0$, and $d' = \beta_1$.

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

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

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

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

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

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

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

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

387 ([https://lindeloev.github.io/shiny-rt/#34_\(shifted\)_log-normal](https://lindeloev.github.io/shiny-rt/#34_(shifted)_log-normal)). This distribution has
 388 two parameters: μ , σ . μ is the mean of the logNormal distribution, and σ is the disperse of
 389 the distribution. The log normal distribution can be extended to shifted log normal
 390 distribution, with one more parameter: shift, which is the earliest possible response.

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

391 Shifted log-normal distribution:

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

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

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

$$\sigma_j \sim Cauchy()$$

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

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

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

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

404 **Synthesized results.** We also reported the synthesized results from the
405 experiments, because many of them shared the similar experimental design. We reported
406 the results in five parts: valence effect, explicit interaction between valence and
407 self-relevance, implicit interaction between valence and self-relevance, specificity of valence
408 effect, and behavior-questionnaire correlation.

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

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

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

417 M_1 is the mean of the first condition, sd_1 is the standard deviation of the first
418 condition, while M_2 is the mean of the second condition, sd_2 is the standard deviation of
419 the second condition. r is the correlation coefficient between data from first and second
420 condition. n is the number of data point (in our case the number of participants included
421 in our research).

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

⁴²⁶ mentioned above, 24 participants were intentionally recruited to participate both exp 1a
⁴²⁷ and exp 2, we only included their results from experiment 1a in the meta-analysis.

⁴²⁸ We also estimated the synthesized effect size using Bayesian hierarchical model,
⁴²⁹ which extended the two-level hierarchical model in each experiment into three-level model,
⁴³⁰ which experiment as an additional level. For SDT, we can use a nested hierarchical model
⁴³¹ to model all the experiment with similar design:

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

⁴³² where

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

⁴³³ The outcomes y_{ijk} are 0 if participant j in experiment k responded “nonmatch” on trial i ,
⁴³⁴ 1 if they responded “match”.

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

⁴³⁵ and the experiment level parameter mu_{0k} and mu_{1k} is from a higher order
⁴³⁶ distribution:

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

⁴³⁷ in which μ_0 and μ_1 means the population level parameter.

⁴³⁸ This model can be easily expand to three-level model in which participants and
⁴³⁹ experiments are two group level variable and participants were nested in the experiments.

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

⁴⁴⁰ y_{ijk} is the RT of the i th trial of the j th participants in the k th experiment.

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

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

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

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

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

Valence effect.

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

Explicit interaction between Valence and self-relevance.

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

Implicit interaction between valence and self-relevance.

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

457 4a) or self-referential (experiment 4a) was presented as task-irrelevant stimuli, whether
458 they would facilitate self-referential or valence effect on perceptual decision-making. For
459 experiment 4a, we reported the comparisons between different valence conditions under the
460 self-referential task and other-referential task. For experiment 4b, we first calculated the
461 effect of valence for both self- and other-referential conditions and then compared the effect
462 size of these three contrast from self-referential condition and from other-referential
463 condition. Note that the results were also analyzed in a standard repeated measure
464 ANOVA (see supplementary materials).

465 ***Specificity of the valence effect.***

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

470 ***Behavior-Questionnaire correlation.***

471 Finally, we explored correlation between results from behavioral results and
472 self-reported measures.

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

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

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

489 Then we confirmed the significant results from the data without self- and
490 other-referential (exp1a, 1b, 1c, 2, 5, 6a). In this confirmatory analysis, we used $\alpha =$
491 0.05 and used bootstrap by BootES package (Kirby & Gerlanc, 2013) to estimate the
492 correlation. To avoid false positive, we further determined the threshold for significant by
493 permutation. More specifically, for each pairs that initially with $p < .05$, we randomly
494 shuffle the participants data of each score and calculated the correlation between the
495 shuffled vectors. After repeating this procedure for 5000 times, we choose arrange these
496 5000 correlation coefficients and use the 95% percentile number as our threshold.

497 **Part 1: Moral valence effect**

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

501 **Experiment 1a**

502 **Methods.**

503 ***Participants.***

504 57 college students (38 female, age = 20.75 ± 2.54 years) participated. 39 of them
505 were recruited from Tsinghua University community in 2014; 18 were recruited from
506 Wenzhou University in 2017. All participants were right-handed except one, and all had

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

511 ***Stimuli and Tasks.***

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

515 ***Procedure.***

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

533 was presented to remind participants to accelerate. Participants were informed of their
534 overall accuracy at the end of each block. The practice phase finished and the experimental
535 task began after the overall performance of accuracy during practice phase achieved 60%.
536 For participants from the Tsinghua community, they completed 6 experimental blocks of 60
537 trials. Thus, there were 60 trials in each condition (bad-person match, bad-person
538 nonmatch, good-person match, good-person nonmatch, neutral-person match, and
539 neutral-person nonmatch). For the participants from Wenzhou University, they finished 6
540 blocks of 120 trials, therefore, 120 trials for each condition.

541 ***Data analysis.***

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

545 **Results.**

546 ***Classic NHST.***

547 *d prime.*

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

550 We found the effect of Valence ($F(1.96, 97.84) = 6.19$, $MSE = 0.27$, $p = .003$,
551 $\hat{\eta}_G^2 = .020$). The post-hoc comparison with multiple comparison correction revealed that
552 the shapes associated with Good-person (2.11, SE = 0.14) has greater *d prime* than shapes
553 associated with Bad-person (1.75, SE = 0.14), $t(50) = 3.304$, $p = 0.0049$. The Good-person
554 condition was also greater than the Neutral-person condition (1.95, SE = 0.16), but didn't
555 reach statistical significant, $t(50) = 1.54$, $p = 0.28$. Neither the Neutral-person condition is
556 significantly greater than the Bad-person condition, $t(50) = 2.109$, $p = .098$.

557 ***Reaction times.***

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

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

560 $MSE = 948.92, p < .001, \hat{\eta}_G^2 = .104$), main effect of valence ($F(1.87, 93.31) = 9.62$,

561 $MSE = 1,673.86, p < .001, \hat{\eta}_G^2 = .016$), and interaction between Matchness and Valence

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

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

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

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

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

567 -2.265, $p = 0.0702$) and Bad Person (728 ms, SE = 11.7), $t(50) = -4.41, p = 0.0002$), and

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

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

570 ***Bayesian hierarchical GLM.***

571 d' prime.

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

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

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

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

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

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

578 95% CI[1.95 2.49]), $P_{PosteriorComparison} = 0.97$. Also, the shapes tagged with neutral

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

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

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

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

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

584 evidence for the difference between good and bad conditions.

585 *Reaction times.*

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

594 **HDDM.**

595 We fitted our data with HDDM, using the response-coding (See also, Hu et al., 2020).
596 We estimated separate drift rate (v), non-decision time (T_0), and boundary separation (a)
597 for each condition. We found that the shapes tagged with good person has higher drift rate
598 and higher boundary separation than shapes tagged with both neutral and bad person.
599 Also, the shapes tagged with neutral person has a higher drift rate than shapes tagged
600 with bad person, but not for the boundary separation. Finally, we found that shapes
601 tagged with bad person had longer non-decision time (see Figure 3).

602 **Experiment 1b**

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

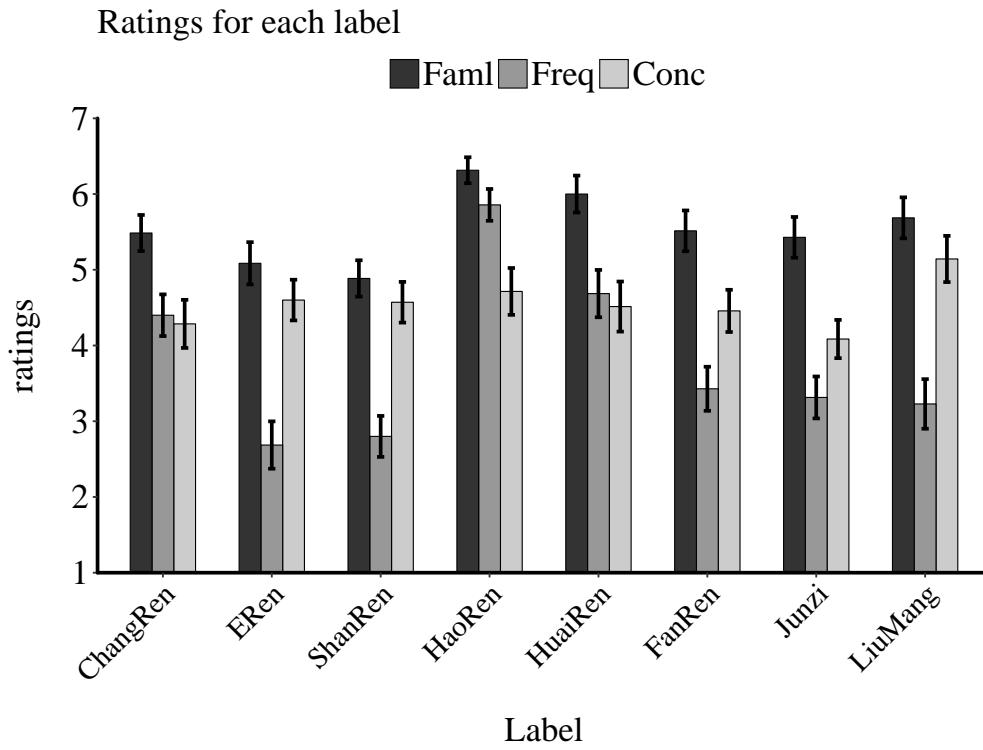
605 **Method.**

606 *Participants.*

607 72 college students (49 female, age = 20.17 ± 2.08 years) participated. 39 of them
608 were recruited from Tsinghua University community in 2014; 33 were recruited from

609 Wenzhou University in 2017. All participants were right-handed except one, and all had
610 normal or corrected-to-normal vision. Informed consent was obtained from all participants
611 prior to the experiment according to procedures approved by the local ethics committees.
612 20 participant's data were excluded from analysis because nearly random level of accuracy,
613 leaving 52 participants (36 female, age = 20.25 ± 2.31 years).

614 **Stimuli and Tasks.** Three geometric shapes (triangle, square, and circle, with 3.7°
615 $\times 3.7^\circ$ of visual angle) were presented above a white fixation cross subtending $0.8^\circ \times 0.8^\circ$
616 of visual angle at the center of the screen. The three shapes were randomly assigned to
617 three labels with different moral valence: a morally bad person (" ", ERen), a morally
618 good person (" ", ShanRen) or a morally neutral person (" ", ChangRen). The order of
619 the associations between shapes and labels was counterbalanced across participants. Three
620 labels used in this experiment is selected based on the rating results from an independent
621 survey, in which participants rated the familiarity, frequency, and concreteness of eight
622 different words online. Of the eight words, three of them are morally positive (HaoRen,
623 ShanRen, Junzi), two of them are morally neutral (ChangRen, FanRen), and three of them
624 are morally negative (HuaiRen, ERen, LiuMang). An independent sample consist of 35
625 participants (22 females, age 20.6 ± 3.11) were recruited to rate these words. Based on the
626 ratings (see supplementary materials Figure S1), we selected ShanRen, ChangRen, and
627 ERen to represent morally positive, neutral, and negative person.



628

Procedure.

629 For participants from both Tsinghua community and Wenzhou community, the
 630 procedure in the current study was exactly same as in experiment 1a.
 631

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

633 **Results.**

634 **NHST.**

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

636 d prime.

637 Repeated measures ANOVA revealed main effect of valence, $F(1.83, 93.20) = 14.98$,
 638 $MSE = 0.18$, $p < .001$, $\hat{\eta}_G^2 = .053$. Paired t test showed that the Good-Person condition
 639 (1.87 ± 0.102) was with greater d prime than Neutral condition (1.44 ± 0.101 , $t(51) =$
 640 5.945 , $p < 0.001$). We also found that the Bad-Person condition (1.67 ± 0.11) has also
 641 greater d prime than neutral condition , $t(51) = 3.132$, $p = 0.008$). There Good-person

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

643 *Reaction times.*

644 We found interaction between Matchness and Valence ($F(1.95, 99.31) = 19.71$,

645 $MSE = 960.92, p < .001, \hat{\eta}_G^2 = .031$) and then analyzed the matched trials and

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

647 of valence $F(1.94, 99.10) = 33.97, MSE = 1,343.19, p < .001, \hat{\eta}_G^2 = .115$. Post-hoc t -tests

648 revealed that shapes associated with Good Person (684 ± 8.77) were responded faster than

649 Neutral-Person (740 ± 9.84), ($t(51) = -8.167, p < 0.001$) and Bad Person (728 ± 9.15),

650 $t(51) = -5.724, p < 0.0001$). While there was no significant differences between Neutral and

651 Bad-Person condition ($t(51) = 1.686, p = 0.221$). For non-matched trials, there was no

652 significant effect of Valence ($F(1.90, 97.13) = 1.80, MSE = 430.15, p = .173, \hat{\eta}_G^2 = .003$).

653 **BGLM.**

654 *Signal detection theory analysis of accuracy.*

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

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

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

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

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

660 greater than shapes tagged with neutral person (2.23, 95% CI[1.95 2.49]),

661 $P_{PosteriorComparison} = 0.97$. Also, the shapes tagged with neutral person is greater than

662 shapes tagged with morally bad person, $P_{PosteriorComparison} = 0.92$.

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

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

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

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

667 evidence for the difference between good and bad conditions.

668 *Reaction time.*

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

677 **HDDM.**

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

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

688 **Experiment 1c**

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

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

693 Method.**694 *Participants.***

695 23 college students (15 female, age = 22.61 ± 2.62 years) participated. All of them
696 were recruited from Tsinghua University community in 2014. Informed consent was
697 obtained from all participants prior to the experiment according to procedures approved by
698 the local ethics committees. No participant was excluded because they overall accuracy
699 were above 0.6.

700 *Stimuli and Tasks.*

701 Three geometric shapes (triangle, square, and circle, with $3.7^\circ \times 3.7^\circ$ of visual angle)
702 were presented above a white fixation cross subtending $0.8^\circ \times 0.8^\circ$ of visual angle at the
703 center of the screen. The three most common names were chosen, which are neutral in
704 moral valence before the manipulation. Three names (Zhang, Wang, Li) were first paired
705 with three paragraphs of behavioral description. Each description includes one sentence of
706 biographic information and four sentences that describing the moral behavioral under that
707 name. To assess the that these three descriptions represented good, neutral, and bad
708 valence, we collected the ratings of three person on six dimensions: morality, likability,
709 trustworthiness, dominance, competence, and aggressiveness, from an independent sample
710 ($n = 34$, 18 female, age = 19.6 ± 2.05). The rating results showed that the person with
711 morally good behavioral description has higher score on morality ($M = 3.59$, $SD = 0.66$)
712 than neutral ($M = 0.88$, $SD = 1.1$), $t(33) = 12.94$, $p < .001$, and bad conditions ($M = -3.4$,
713 $SD = 1.1$), $t(33) = 30.78$, $p < .001$. Neutral condition was also significant higher than bad
714 conditions $t(33) = 13.9$, $p < .001$ (See supplementary materials).

715 *Procedure.*

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

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

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

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

Signal detection theory analysis of accuracy.

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

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

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

748 ***Reaction time.***

749 We fitted a Bayesian hierarchical GLM for RTs, with a log-normal distribution as the
750 link function. We used the posterior distribution of the regression coefficient to make
751 statistical inferences. As in previous studies, the matched conditions are much faster than
752 the mismatched trials ($P_{PosteriorComparison} = .75$). We focused on matched trials only, and
753 compared different conditions: Good () is not faster than the neutral (),
754 $P_{PosteriorComparison} = .5$, it was faster than the Bad condition (),
755 $P_{PosteriorComparison} = .88$. And the neutral condition is faster than the bad condition,
756 $P_{PosteriorComparison} = .95$. However, the mismatched trials are largely overlapped.

757 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et
758 al., 2020). We estimated separate drift rate (v), non-decision time (T_0), and boundary
759 separation (a) for each condition. We found that the shapes tagged with good person has
760 higher drift rate and higher boundary separation than shapes tagged with both neutral and
761 bad person. Also, the shapes tagged with neutral person has a higher drift rate than
762 shapes tagged with bad person, but not for the boundary separation. Finally, we found
763 that shapes tagged with bad person had longer non-decision time (see figure 9)).

764 **Experiment 2: Sequential presenting**

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

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

777 **Method.**

778 ***Participants.***

779 35 participants (17 female, age = 21.66 ± 3.03) were recruited. 24 of them had
780 participated in Experiment 1a (9 male, mean age = 21.9, s.d. = 2.9), and the time gap
781 between these experiment 1a and experiment 2 is at least six weeks. The results of 1
782 participants were excluded from analysis because of less than 60% overall accuracy,
783 remains 34 participants (17 female, age = 21.74 ± 3.04).

784 ***Procedure.***

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

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

798 ***Data analysis.***

799 Data was analyzed as in study 1a.

800 **Results.**

801 **NHST.**

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

804 *d prime.*

805 There was evidence for the main effect of valence, $F(1.83, 60.36) = 14.41$,
806 $MSE = 0.23$, $p < .001$, $\hat{\eta}_G^2 = .066$. Paired t test showed that the Good-Person condition
807 (2.79 ± 0.17) was with greater d prime than Netural condition (2.21 ± 0.16 , $t(33) = 4.723$,
808 $p = 0.001$) and Bad-person condition (2.41 ± 0.14), $t(33) = 4.067$, $p = 0.008$). There was
809 no-significant difference between Neutral-person and Bad-person conidition, $t(33) = -1.802$,
810 $p = 0.185$.

811 *Reaction time.*

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

814 We found interaction between Matchness and Valence ($F(1.99, 65.70) = 9.53$,
815 $MSE = 605.36$, $p < .001$, $\hat{\eta}_G^2 = .017$) and then analyzed the matched trials and
816 mismatched trials separately, as in experiment 1a. For matched trials, we found the effect
817 of valence $F(1.99, 65.76) = 10.57$, $MSE = 1,192.65$, $p < .001$, $\hat{\eta}_G^2 = .067$. Post-hoc t -tests
818 revealed that shapes associated with Good Person (548 ± 9.4) were responded faster than
819 Neutral-Person (582 ± 10.9), ($t(33) = -3.95$, $p = 0.0011$) and Bad Person (582 ± 10.2),
820 $t(33) = -3.9$, $p = 0.0013$). While there was no significant differences between Neutral and

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

823 **BGLMM.**

824 *Signal detection theory analysis of accuracy.*

825 We fitted a Bayesian hierarchical GLM for SDT. The results showed that when the
 826 shapes were tagged with labels with different moral valence, the sensitivity (d') and criteria
 827 (c) were both influence. For the d' , we found that the shapes tagged with morally good
 828 person (2.46, 95% CI[2.21 2.72]) is greater than shapes tagged with moral bad (2.07, 95%
 829 CI[1.83 2.32]), $P_{PosteriorComparison} = 1$. Shape tagged with morally good person is also
 830 greater than shapes tagged with neutral person (2.23, 95% CI[1.95 2.49]),
 831 $P_{PosteriorComparison} = 0.97$. Also, the shapes tagged with neutral person is greater than
 832 shapes tagged with morally bad person, $P_{PosteriorComparison} = 0.92$.

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

838 *Reaction times.*

839 We fitted a Bayesian hierarchical GLM for RTs, with a log-normal distribution as the
 840 link function. We used the posterior distribution of the regression coefficient to make
 841 statistical inferences. As in previous studies, the matched conditions are much faster than
 842 the mismatched trials ($P_{PosteriorComparison} = .75$). We focused on matched trials only, and
 843 compared different conditions: Good () is not faster than the neutral (),
 844 $P_{PosteriorComparison} = .5$, it was faster than the Bad condition (),
 845 $P_{PosteriorComparison} = .88$. And the neutral condition is faster than the bad condition,
 846 $P_{PosteriorComparison} = .95$. However, the mismatched trials are largely overlapped.

847 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et
848 al., 2020). We estimated separate drift rate (v), non-decision time (T_0), and boundary
849 separation (a) for each condition. We found that the shapes tagged with good person has
850 higher drift rate and higher boundary separation than shapes tagged with both neutral and
851 bad person. Also, the shapes tagged with neutral person has a higher drift rate than
852 shapes tagged with bad person, but not for the boundary separation. Finally, we found
853 that shapes tagged with bad person had longer non-decision time (see figure
854 @ref(fig:plot-exp1c -HDDM))).

855 Discussion

856 In this experiment, we repeated the results pattern that the positive moral valenced
857 stimuli has an advantage over the neutral or the negative valence association. Moreover,
858 with a cross-task analysis, we did not find evidence that the experiment task interacted
859 with moral valence, suggesting that the effect might not be effect by experiment task.

860 These findings suggested that the facilitation effect of positive moral valence is robust and
861 not affected by task. This robust effect detected by the associative learning is unexpected.

862 Experiment 6a: EEG study 1

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

865 Method.

866 *Participants.*

867 24 college students (8 female, age = 22.88 ± 2.79) participated the current study, all
868 of them were from Tsinghua University in 2014. Informed consent was obtained from all
869 participants prior to the experiment according to procedures approved by a local ethics
870 committee. No participant was excluded from behavioral analysis.

871 **Experimental design.** The experimental design of this experiment is same as
872 experiment 2: a 3×2 within-subject design with moral valence (good, neutral and bad
873 associations) and matchness between shape and label (match vs. mismatch for the personal
874 association) as within-subject variables.

875 *Stimuli.*

876 Three geometric shapes (triangle, square and circle, each $4.6^\circ \times 4.6^\circ$ of visual angle)
877 were presented at the center of screen for 50 ms after 500ms of fixation ($0.8^\circ \times 0.8^\circ$ of
878 visual angle). The association of the three shapes to bad person (“ , HuaiRen”), good
879 person (“ , HaoRen”) or ordinary person (“ , ChangRen”) was counterbalanced across
880 participants. The words bad person, good person or ordinary person ($3.6^\circ \times 1.6^\circ$) was also
881 displayed at the center fo the screen. Participants had to judge whether the pairings of
882 label and shape matched (e.g., Does the circle represent a bad person?). The experiment
883 was run on a PC using E-prime software (version 2.0). These stimuli were displayed on a
884 22-in CRT monitor (1024×768 at 100Hz). We used backward masking to avoid
885 over-processing of the moral words, in which a scrambled picture were presented for 900 ms
886 after the label. Also, to avoid the ceiling effect on accuracy, shapes were presented on a
887 noisy background based on our pilot studies. The noisy images were made by scrambling a
888 picture of 3/4gray and 1/4 white at resolution of 2×2 pixel.

889 *Procedure.*

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

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

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

901 All the stimuli were presented on a gray background (RGB: 127, 127, 127). E-primed
902 2.0 was used to present stimuli and collect behavioral results. Data were collected and
903 analyzed when accuracy performance in total reached 60%.

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

905 **Results.**

906 **NHST.**

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

909 *d prime.*

910 We conducted repeated measures ANOVA, with moral valence as independent
911 variable. The results revealed the main effect of valence ($F(1.74, 40.05) = 3.76$,
912 $MSE = 0.10$, $p = .037$, $\eta^2_G = .021$). Post-hoc analysis revealed that shapes link with Good
913 person (mean = 3.13, SE = 0.109) is greater than Neutral condition (mean = 2.88, SE =
914 0.14), $t = 2.916$, $df = 24$, $p = 0.02$, p-value adjusted by Tukey method, but the *d* prime
915 between Good and bad (mean = 3.03, SE = 0.142) ($t = 1.512$, $df = 24$, $p = 0.3034$, p-value
916 adjusted by Tukey method), bad and neutral ($t = 1.599$, $df = 24$, $p = 0.2655$, p-value
917 adjusted by Tukey method) were not significant.

918 *Reaction times.*

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

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

922 $MSE = 450.47, p < .001, \hat{\eta}_G^2 = .021$) and then analyzed the matched trials and
 923 mismatched trials separately, as in experiment 2. For matched trials, we found the effect of
 924 valence $F(1.97, 45.25) = 32.37, MSE = 522.42, p < .001, \hat{\eta}_G^2 = .078$. For non-matched
 925 trials, there was no significant effect of Valence ($F(1.77, 40.67) = 0.35, MSE = 242.15,$
 926 $p = .679, \hat{\eta}_G^2 = .000$). Post-hoc t -tests revealed that shapes associated with Good Person
 927 (mean = 550, SE = 13.8) were responded faster than Neutral-Person (501, SE = 14.7),
 928 ($t(24) = -5.171, p = 0.0001$) and Bad Person (523, SE = 16.3), $t(24) = -8.137, p <$
 929 0.0001),, and Neutral is faster than Bad-Person condition ($t(32) = -3.282, p = 0.0085$).

930 **BGLM.**

931 *Signal detection theory analysis of accuracy.*

932 *Reaction time.*

933 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et
 934 al., 2020). We estimated separate drift rate (v), non-decision time (T_0), and boundary
 935 separation (a) for each condition. We found that, similar to experiment 2, the shapes
 936 tagged with good person has higher drift rate and higher boundary separation than shapes
 937 tagged with both neutral and bad person, but only for the self-referential condition. Also,
 938 the shapes tagged with neutral person has a higher drift rate than shapes tagged with bad
 939 person, but not for the boundary separation, and this effect also exist only for the
 940 self-referential condition.

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

Part 2: interaction between valence and identity

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

949 Experiment 3a

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

954 Method.**955 Participants.**

956 38 college students (15 female, age = 21.92 ± 2.16) participated in experiment 3a.
957 All of them were right-handed, and all had normal or corrected-to-normal vision. Informed
958 consent was obtained from all participants prior to the experiment according to procedures
959 approved by a local ethics committee. One female and one male student did not finish the
960 experiment, and 1 participants' data were excluded from analysis because less than 60%
961 overall accuracy, remains 35 participants (13 female, age = 22.11 ± 2.13).

962 Design.

963 Study 3a combined moral valence with self-relevance, hence the experiment has a $2 \times$
964 3×2 within-subject design. The first variable was self-relevance, include two levels:
965 self-relevance vs. stranger-relevance; the second variable was moral valence, include good,
966 neutral and bad; the third variable was the matching between shape and label: match
967 vs. nonmatch.

968 Stimuli.

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

974 Procedure.

975 After being fully explained and signed the informed consent, participants were
976 instructed to chose a name that can represent a stranger with same gender as the
977 participant themselves, from a common Chinese name pool. Before experiment, the
978 experimenter explained the meaning of each label to participants. For example, the “good
979 self” mean the morally good side of themselves, them could imagine the moment when they
980 do something’s morally applauded, “bad self” means the morally bad side of themselves,
981 they could also imagine the moment when they doing something morally wrong, and
982 “neutral self” means the aspect of self that does not related to morality, they could imagine
983 the moment when they doing something irrelevant to morality. In the same sense, the
984 “good other”, “bad other”, and “neutral other” means the three different aspects of the
985 stranger, whose name was chosen before the experiment. Then, the experiment proceeded
986 as study 1a. Each participant finished 6 blocks, each have 120 trials. The sequence of trials
987 was pseudo-randomized so that there are 10 matched trials for each condition and 10
988 non-matched trials for each condition (good self, neutral self, bad self, good other, neutral
989 other, bad other) for each block.

990 Data Analysis.

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

994 **Results.**

995 **NHST.**

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

998 *d prime.*

999 There was evidence for the main effect of valence, $F(1.89, 64.37) = 11.09$,

1000 $MSE = 0.23$, $p < .001$, $\hat{\eta}_G^2 = .039$, and main effect of self-relevance, $F(1, 34) = 3.22$,

1001 $MSE = 0.54$, $p = .082$, $\hat{\eta}_G^2 = .015$, as well as the interaction, $F(1.79, 60.79) = 3.39$,

1002 $MSE = 0.43$, $p = .045$, $\hat{\eta}_G^2 = .022$.

1003 We then conducted separated ANOVA for self-referential and other-referential trials.

1004 The valence effect was shown for the self-referential conditions, $F(1.65, 56.25) = 13.98$,

1005 $MSE = 0.31$, $p < .001$, $\hat{\eta}_G^2 = .119$. Post-hoc test revealed that the Good-Self condition

1006 (1.97 ± 0.14) was with greater d prime than Neutral condition $(1.41 \pm 0.12$, $t(34) = 4.505$,

1007 $p = 0.0002$), and Bad-self condition (1.43 ± 0.102) , $t(34) = 3.856$, $p = 0.0014$. There was

1008 difference between neutral and bad condition, $t(34) = -0.238$, $p = 0.9694$. However, no

1009 effect of valence was found for the other-referential condition $F(1.98, 67.36) = 0.38$,

1010 $MSE = 0.35$, $p = .681$, $\hat{\eta}_G^2 = .004$.

1011 *Reaction time.*

1012 We found interaction between Matchness and Valence ($F(1.98, 67.44) = 26.29$,

1013 $MSE = 730.09$, $p < .001$, $\hat{\eta}_G^2 = .025$) and then analyzed the matched trials and nonmatch

1014 trials separately, as in previous experiments.

1015 For the match trials, we found that the interaction between identity and valence,

1016 $F(1.72, 58.61) = 3.89$, $MSE = 2,750.19$, $p = .032$, $\hat{\eta}_G^2 = .019$, as well as the main effect of

1017 valence $F(1.98, 67.34) = 35.76$, $MSE = 1,127.25$, $p < .001$, $\hat{\eta}_G^2 = .079$, but not the effect of

1018 identity $F(1, 34) = 0.20$, $MSE = 3,507.14$, $p = .660$, $\hat{\eta}_G^2 = .001$. As for the d prime, we

1019 separated analyzed the self-referential and other-referential trials. For the Self-referential
 1020 trials, we found the main effect of valence, $F(1.80, 61.09) = 30.39$, $MSE = 1,584.53$,
 1021 $p < .001$, $\hat{\eta}_G^2 = .159$; for the other-referential trials, the effect of valence is weaker,
 1022 $F(1.86, 63.08) = 2.85$, $MSE = 2,224.30$, $p = .069$, $\hat{\eta}_G^2 = .024$. We then focused on the self
 1023 conditions: the good-self condition (713 ± 12) is faster than neutral- (776 ± 11.8), $t(34) =$
 1024 -7.396 , $p < .0001$, and bad-self (772 ± 10.1) conditions, $t(34) = -5.66$, $p < .0001$. But
 1025 there is not difference between neutral- and bad-self conditions, $t(34) = 0.481$, $p = 0.881$.

1026 For the nonmatch trials, we didn't found any strong effect: identity, $F(1, 34) = 3.43$,
 1027 $MSE = 660.02$, $p = .073$, $\hat{\eta}_G^2 = .004$, valence $F(1.89, 64.33) = 0.40$, $MSE = 444.10$,
 1028 $p = .661$, $\hat{\eta}_G^2 = .001$, or interaction between the two $F(1.94, 66.02) = 2.42$, $MSE = 817.35$,
 1029 $p = .099$, $\hat{\eta}_G^2 = .007$.

1030 **BGLM.**

1031 *Signal detection theory analysis of accuracy.*

1032 We found that the d prime is greater when shapes were associated with good self
 1033 condition than with neutral self or bad self, but shapes associated with bad self and neutral
 1034 self didn't show differences. Comparing the self vs other under three condition revealed
 1035 that shapes associated with good self is greater than with good other, but with a weak
 1036 evidence. In contrast, for both neutral and bad valence condition, shapes associated with
 1037 other had greater d prime than with self.

1038 *Reaction time.*

1039 In reaction times, we found that same trends in the match trials as in the RT: while
 1040 the shapes associated with good self was greater than with good other (log mean diff =
 1041 -0.02858 , 95%HPD[-0.070898, 0.0154]), the direction is reversed for neutral and negative
 1042 condition. see Figure 17

1043 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et
 1044 al., 2020). We estimated separate drift rate (v), non-decision time (T_0), and boundary

1045 separation (*a*) for each condition. We found that the shapes tagged with good person has
1046 higher drift rate and higher boundary separation than shapes tagged with both neutral and
1047 bad person. Also, the shapes tagged with neutral person has a higher drift rate than
1048 shapes tagged with bad person, but not for the boundary separation. Finally, we found
1049 that shapes tagged with bad person had longer non-decision time (see figure 18)).

1050 **Experiment 3b**

1051 In study 3a, participants had to remember 6 pairs of association, which cause high
1052 cognitive load during the whole experiment. To eliminate the influence of cognitive load, we
1053 conducted study 3b, in which participant learn three aspect of self and stranger separately
1054 in to consecutive task. We hypothesize that we will replicate the pattern of study 3a, i.e.,
1055 the effect of moral valence only occurs for self-relevant conditions. #### Method

1056 **Participants.**

1057 Study 3b were finished in 2017, at that time we have calculated that the effect size
1058 (Cohen's *d*) of good-person (or good-self) vs. bad-person (or bad-other) was between $0.47 \sim 0.53$, based on the data from Tsinghua community in study 1a, 1b, 2, 3a, 4a, and 4b. Based
1059 on this effect size, we estimated that 54 participants would allow we to detect the effect
1060 size of Cohen's $= 0.5$ with 95% power and alpha = 0.05, using G*power 3.192 (Faul,
1061 Erdfelder, Buchner, & Lang, 2009). Therefore, we planned to stop after we arrived this
1062 number. During the data collected at Wenzhou University, 61 participants (45 females; 19
1063 to 25 years of age, age = 20.42 ± 1.77) came to the testing room and we tested all of them
1064 during a single day. All participants were right-handed, and all had normal or
1065 corrected-to-normal vision. Informed consent was obtained from all participants prior to
1066 the experiment according to procedures approved by a local ethics committee. 4
1067 participants' data were excluded from analysis because their over all accuracy was lower
1068 than 60%, 1 more participant was excluded because of zero hit rate for one condition,
1069 leaving 56 participants (43 females; 19 to 25 years old, age = 20.27 ± 1.60).

Design.

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

Stimuli.

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

Procedure.

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

1097 ***Data Analysis.***

1098 Same as experiment 3a.

1099 **Results.**

1100 ***NHST.***

1101 Figure 19 shows d prime and reaction times of experiment 3b. Less than 5% correct
 1102 trials with less than 200ms reaction times were excluded.

1103 *d prime.*

1104 There was no evidence for the main effect of valence, $F(1.92, 105.43) = 1.90$,

1105 $MSE = 0.33$, $p = .157$, $\hat{\eta}_G^2 = .005$, but we found a main effect of self-relevance,

1106 $F(1, 55) = 4.65$, $MSE = 0.89$, $p = .035$, $\hat{\eta}_G^2 = .017$, as well as the interaction,

1107 $F(1.90, 104.36) = 5.58$, $MSE = 0.26$, $p = .006$, $\hat{\eta}_G^2 = .011$.

1108 We then conducted separated ANOVA for self-referential and other-referential trials.

1109 The valence effect was shown for the self-referential conditions, $F(1.75, 96.42) = 6.73$,

1110 $MSE = 0.30$, $p = .003$, $\hat{\eta}_G^2 = .037$. Post-hoc test revealed that the Good-Self condition

1111 (2.15 ± 0.12) was with greater d prime than Neutral condition $(1.83 \pm 0.12$, $t(34) = 3.36$,

1112 $p = 0.0031$), and Bad-self condition (1.87 ± 0.12) , $t(34) = 2.955$, $p = 0.01$. There was

1113 difference between neutral and bad condition, $t(34) = -0.039$, $p = 0.914$. However, no effect

1114 of valence was found for the other-referential condition $F(1.93, 105.97) = 0.61$,

1115 $MSE = 0.31$, $p = .539$, $\hat{\eta}_G^2 = .002$.

1116 *Reaction time.*

1117 We found interaction between Matchness and Valence ($F(1.86, 102.47) = 15.44$,

1118 $MSE = 3, 112.78$, $p < .001$, $\hat{\eta}_G^2 = .006$) and then analyzed the matched trials and

1119 nonmatch trials separately, as in previous experiments.

1120 For the match trials, we found that the interaction between identity and valence,

1121 $F(1.67, 92.11) = 6.14$, $MSE = 6, 472.48$, $p = .005$, $\hat{\eta}_G^2 = .009$, as well as the main effect of

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

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

BGLM.

Signal detection theory analysis of accuracy.

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

Reaction time.

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

1148 condition. see Figure 20

1149 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et
1150 al., 2020). We estimated separate drift rate (v), non-decision time (T_0), and boundary
1151 separation (a) for each condition. We found that, similar to experiment 3a, the shapes
1152 tagged with good person has higher drift rate and higher boundary separation than shapes
1153 tagged with both neutral and bad person, but only for the self-referential condition. Also,
1154 the shapes tagged with neutral person has a higher drift rate than shapes tagged with bad
1155 person, but not for the boundary separation, and this effect also exist only for the
1156 self-referential condition.

1157 Interestingly, we found that in both self-referential and other-referential conditions,
1158 the shapes associated bad valence have higher drift rate and higher boundary separation.
1159 which might suggest that the shape associated with bad stimuli might be prioritized in the
1160 non-match trials (see figure 21)).

1161 Experiment 6b

1162 Experiment 6b was conducted to study the neural correlates of the prioritization
1163 effect of positive self, i.e., the neural underlying of the behavioral effect found int
1164 experiment 3a. However, as in experiment 6a, the procedure of this experiment was
1165 modified to adopted to ERP experiment.

1166 Method.

1167 Participants.

1168 23 college students (8 female, age = 22.86 ± 2.47) participated the current study, all
1169 of them were recruited from Tsinghua University in 2016. Informed consent was obtained
1170 from all participants prior to the experiment according to procedures approved by a local
1171 ethics committee. For day 1's data, 1 participant was excluded from the current analysis
1172 because of lower than 60% overall accuracy, remaining 22 participants (8 female, age =

₁₁₇₃ 22.76 ± 2.49). For day 2's data, one participant dropped out, leaving 22 participants (9
₁₁₇₄ female, age = 23.05 ± 2.46), all of them has overall accuracy higher than 60%.

₁₁₇₅ ***Design.***

₁₁₇₆ The experimental design of this experiment is same as experiment 3: a 2 × 3 × 2
₁₁₇₇ within-subject design with self-relevance (self-relevant vs. other-relevant), moral valence
₁₁₇₈ (good, neutral, and bad) and matchness between shape and label (match vs. mismatch) as
₁₁₇₉ within-subject variables.

₁₁₈₀ ***Stimuli.***

₁₁₈₁ As in experiment 3a, 6 shapes were included (triangle, square, circle, trapezoid,
₁₁₈₂ diamond, regular pentagon), as well as 6 labels (good self, neutral self, bad self, good
₁₁₈₃ person, bad person, neutral person). To match the concreteness of the label, we asked
₁₁₈₄ participant to chosen an unfamiliar name of their own gender to be the stranger.

₁₁₈₅ ***Procedure.***

₁₁₈₆ The procedure was similar to Experiment 2 and 6a. Subjects first learned the
₁₁₈₇ associations between labels and shapes and then completed a shape-label matching task. In
₁₁₈₈ each trial of the matching task, a fixation were first presented for 500 ms, followed by a 50
₁₁₈₉ ms label; then, a scrambled picture presented 900 ms. After the backward mask, the shape
₁₁₉₀ were presented on a noisy background for 50ms. Participant have to response in 1000ms
₁₁₉₁ after the presentation of the shape, and finally, a feedback screen was presented for 500 ms.
₁₁₉₂ The inter-trial interval (ITI) were randomly varied at the range of 1000 ~ 1400 ms.

₁₁₉₃ All the stimuli were presented on a gray background (RGB: 127, 127, 127). E-primed
₁₁₉₄ 2.0 was used to present stimuli and collect behavioral results. Data were collected and
₁₁₉₅ analyzed when accuracy performance in total reached 60%.

₁₁₉₆ Because learning 6 associations was more difficult than 3 associations and participant
₁₁₉₇ might have low accuracy (see experiment 3a), the current study had extended to a two-day

1198 paradigm to maximizing the accurate trials that can be used in EEG data. At the first day,
1199 participants learnt the associations and finished 9 blocks of the matching task, each had
1200 120 trials, without EEG recording. That is, each condition has 90 trials.

1201 Participants came back to lab at the second day and finish the same task again, with
1202 EEG recorded. Before the EEG experiment, each participant finished a practice session
1203 again, if their accuracy is equal or higher than 85%, they start the experiment (one
1204 participant used lower threshold 75%). Each participant finished 18 blocks, each has 90
1205 trials. One participant finished additional 6 blocks because of high error rate at the
1206 beginning, another two participant finished addition 3 blocks because of the technique
1207 failure in recording the EEG data. To increase the number of trials that can be used for
1208 EEG data analysis, matched trials has twice number as mismatched trials, therefore, for
1209 matched trials each participants finished 180 trials for each condition, for mismatched
1210 trials, each conditions has 90 trials.

1211 ***Data Analysis.***

1212 Same as experiment 3a.

1213 **Results of Day 1.**

1214 ***NHST.***

1215 Figure 22 shows d prime and reaction times of experiment 3b. Less than 5% correct
1216 trials with less than 200ms reaction times were excluded.

1217 ***d prime.***

1218 There was no evidence for the main effect of valence, $F(1.91, 40.20) = 11.98$,
1219 $MSE = 0.15$, $p < .001$, $\hat{\eta}_G^2 = .040$, but we found a main effect of self-relevance,
1220 $F(1, 21) = 1.21$, $MSE = 0.20$, $p = .284$, $\hat{\eta}_G^2 = .003$, as well as the interaction,
1221 $F(1.28, 26.90) = 12.88$, $MSE = 0.21$, $p = .001$, $\hat{\eta}_G^2 = .041$.

1222 We then conducted separated ANOVA for self-referential and other-referential trials.

1223 The valence effect was shown for the self-referential conditions, $F(1.73, 36.42) = 29.31$,
 1224 $MSE = 0.14$, $p < .001$, $\hat{\eta}_G^2 = .147$. Post-hoc test revealed that the Good-Self condition
 1225 (2.15 ± 0.12) was with greater d prime than Neutral condition (1.83 ± 0.12 , $t(34) = 3.36$,
 1226 $p = 0.0031$), and Bad-self condition (1.87 ± 0.12), $t(34) = 2.955$, $p = 0.01$. There was
 1227 difference between neutral and bad condition, $t(34) = -0.039$, $p = 0.914$. However, no effect
 1228 of valence was found for the other-referential condition $F(1.75, 36.72) = 0.00$, $MSE = 0.18$,
 1229 $p = .999$, $\hat{\eta}_G^2 = .000$.

1230 *Reaction time.*

1231 We found interaction between Matchness and Valence ($F(1.79, 37.63) = 4.07$,
 1232 $MSE = 704.90$, $p = .029$, $\hat{\eta}_G^2 = .003$) and then analyzed the matched trials and nonmatch
 1233 trials separately, as in previous experiments.

1234 For the match trials, we found that the interaction between identity and valence,
 1235 $F(1.72, 36.16) = 4.55$, $MSE = 1,560.90$, $p = .022$, $\hat{\eta}_G^2 = .015$, as well as the main effect of
 1236 valence $F(1.93, 40.55) = 9.83$, $MSE = 1,951.84$, $p < .001$, $\hat{\eta}_G^2 = .044$, but not the effect of
 1237 identity $F(1, 21) = 4.87$, $MSE = 2,032.05$, $p = .039$, $\hat{\eta}_G^2 = .012$. As for the d prime, we
 1238 separated analyzed the self-referential and other-referential trials. For the Self-referential
 1239 trials, we found the main effect of valence, $F(1.92, 40.38) = 14.48$, $MSE = 1,647.20$,
 1240 $p < .001$, $\hat{\eta}_G^2 = .112$; for the other-referential trials, the effect of valence is weaker,
 1241 $F(1.79, 37.50) = 1.04$, $MSE = 1,842.07$, $p = .356$, $\hat{\eta}_G^2 = .008$. We then focused on the self
 1242 conditions: the good-self condition (713 ± 12) is faster than neutral- (776 ± 11.8), $t(34) =$
 1243 -7.396 , $p < .0001$, and bad-self (772 ± 10.1) conditions, $t(34) = -5.66$, $p < .0001$. But
 1244 there is not difference between neutral- and bad-self conditions, $t(34) = 0.481$, $p = 0.881$.

1245 For the nonmatch trials, we didn't found any strong effect: identity, $F(1, 21) = 2.76$,
 1246 $MSE = 1,718.93$, $p = .112$, $\hat{\eta}_G^2 = .006$, valence $F(1.61, 33.77) = 3.81$, $MSE = 1,532.21$,
 1247 $p = .041$, $\hat{\eta}_G^2 = .012$, or interaction between the two $F(1.90, 39.97) = 2.23$, $MSE = 720.80$,
 1248 $p = .123$, $\hat{\eta}_G^2 = .004$.

BGLM.*Signal detection theory analysis of accuracy.*

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

Reaction time.

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

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

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

1274

Part 3: Implicit binding between valence and identity

1275

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

1278

Experiment 4a: Morality as task-irrelevant variable

1279

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

1288

Method.

1289

Participants.

1290

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

Design.

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

Stimuli.

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

Procedure.

The procedure was similar to Experiment 1. There were 6 blocks of trial, each with

₁₃₂₄ 120 trials for 2017 data. Due to procedure error, the data collected in 2015 in Tsinghua
₁₃₂₅ community only have 60 trials for each block, i.e., 30 trials per condition.

₁₃₂₆ As in study 3a, before each task, the instruction showed the meaning of each label to
₁₃₂₇ participants. The self-matching task and other-matching task were randomized between
₁₃₂₈ participants. Each participant finished 6 blocks, each have 120 trials.

₁₃₂₉ ***Data Analysis.***

₁₃₃₀ Same as experiment 3a.

₁₃₃₁ **Results.**

₁₃₃₂ ***NHST.***

₁₃₃₃ Figure 25 shows d prime and reaction times of experiment 3a. Less than 5% correct
₁₃₃₄ trials with less than 200ms reaction times were excluded.

₁₃₃₅ d prime.

₁₃₃₆ There was no evidence for the main effect of valence, $F(1.93, 111.66) = 0.53$,
₁₃₃₇ $MSE = 0.12$, $p = .581$, $\hat{\eta}_G^2 = .000$, but we found a main effect of self-relevance,
₁₃₃₈ $F(1, 58) = 121.04$, $MSE = 0.48$, $p < .001$, $\hat{\eta}_G^2 = .189$, as well as the interaction,
₁₃₃₉ $F(1.99, 115.20) = 4.12$, $MSE = 0.14$, $p = .019$, $\hat{\eta}_G^2 = .004$.

₁₃₄₀ We then conducted separated ANOVA for self-referential and other-referential trials.

₁₃₄₁ The valence effect was shown for the self-referential conditions, $F(1.95, 112.92) = 3.01$,
₁₃₄₂ $MSE = 0.15$, $p = .055$, $\hat{\eta}_G^2 = .008$. Post-hoc test revealed that the Good-Self condition
₁₃₄₃ (2.15 ± 0.12) was with greater d prime than Neutral condition (1.83 ± 0.12 , $t(34) = 3.36$,
₁₃₄₄ $p = 0.0031$), and Bad-self condition (1.87 ± 0.12), $t(34) = 2.955$, $p = 0.01$. There was
₁₃₄₅ difference between neutral and bad condition, $t(34) = -0.039$, $p = 0.914$. However, no effect
₁₃₄₆ of valence was found for the other-referential condition $F(1.98, 114.61) = 1.75$,
₁₃₄₇ $MSE = 0.10$, $p = .179$, $\hat{\eta}_G^2 = .003$.

₁₃₄₈ Reaction time.

1349 We found interaction between Matchness and Valence ($F(1.94, 112.64) = 0.84$,
 1350 $MSE = 465.35, p = .432, \hat{\eta}_G^2 = .000$) and then analyzed the matched trials and nonmatch
 1351 trials separately, as in previous experiments.

1352 For the match trials, we found that the interaction between identity and valence,
 1353 $F(1.90, 110.18) = 4.41, MSE = 465.91, p = .016, \hat{\eta}_G^2 = .003$, as well as the main effect of
 1354 valence $F(1.98, 114.82) = 0.94, MSE = 606.30, p = .392, \hat{\eta}_G^2 = .001$, but not the effect of
 1355 identity $F(1, 58) = 124.15, MSE = 4,037.53, p < .001, \hat{\eta}_G^2 = .257$. As for the d prime, we
 1356 separated analyzed the self-referential and other-referential trials. For the Self-referential
 1357 trials, we found the main effect of valence, $F(1.97, 114.32) = 6.29, MSE = 367.25$,
 1358 $p = .003, \hat{\eta}_G^2 = .006$; for the other-referential trials, the effect of valence is weaker,
 1359 $F(1.95, 112.89) = 0.35, MSE = 699.50, p = .699, \hat{\eta}_G^2 = .001$. We then focused on the self
 1360 conditions: the good-self condition (713 ± 12) is faster than neutral- (776 ± 11.8), $t(34) =$
 1361 $-7.396, p < .0001$, and bad-self (772 ± 10.1) conditions, $t(34) = -5.66, p < .0001$. But
 1362 there is not difference between neutral- and bad-self conditions, $t(34) = 0.481, p = 0.881$.

1363 For the nonmatch trials, we didn't found any strong effect: identity, $F(1, 58) = 0.16$,
 1364 $MSE = 1,547.37, p = .692, \hat{\eta}_G^2 = .000$, valence $F(1.96, 113.52) = 0.68, MSE = 390.26$,
 1365 $p = .508, \hat{\eta}_G^2 = .000$, or interaction between the two $F(1.90, 110.27) = 0.04$,
 1366 $MSE = 585.80, p = .953, \hat{\eta}_G^2 = .000$.

1367 **BGLM.**

1368 *Signal detection theory analysis of accuracy.*

1369 We found that the d prime is greater when shapes were associated with good self
 1370 condition than with neutral self or bad self, but shapes associated with bad self and neutral
 1371 self didn't show differences. comparing the self vs other under three condition revealed that
 1372 shapes associated with good self is greater than with good other, but with a weak evidence.
 1373 In contrast, for both neutral and bad valence condition, shapes associated with other had
 1374 greater d prime than with self.

1375 *Reaction time.*

1376 In reaction times, we found that same trends in the match trials as in the RT: while
1377 the shapes associated with good self was greater than with good other (log mean diff =
1378 -0.02858, 95%HPD[-0.070898, 0.0154]), the direction is reversed for neutral and negative
1379 condition. see Figure 26

1380 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et
1381 al., 2020). We estimated separate drift rate (v), non-decision time (T_0), and boundary
1382 separation (a) for each condition. We found that the shapes tagged with good person has
1383 higher drift rate and higher boundary separation than shapes tagged with both neutral and
1384 bad person. Also, the shapes tagged with neutral person has a higher drift rate than
1385 shapes tagged with bad person, but not for the boundary separation. Finally, we found
1386 that shapes tagged with bad person had longer non-decision time (see figure 27)).

1387 **Experiment 4b: Morality as task-irrelevant variable**

1388 In study 4b, we changed the role of valence and identity in task. In this experiment,
1389 participants learn the association between moral valence and the made perceptual match
1390 judgments to associations between different moral valence and shapes as in study 1-3.
1391 Different from experiment 1 ~ 3, we made put the labels of “self/other” in the shapes so
1392 that identity served as an task irrelevant variable. As in experiment 4b, we also
1393 hypothesized that the intrinsic binding between morally good and self will enhance the
1394 performance of good self condition, even identity is irrelevant to the task.

1395 **Method.**

1396 **Participants.**

1397 53 participants (39 female, age = 20.57 ± 1.81) participated the current study, 34 of
1398 them were from Tsinghua University in 2015, 19 were from Wenzhou University
1399 participated in 2017. All participants were right-handed, and all had normal or

1400 corrected-to-normal vision. Informed consent was obtained from all participants prior to
1401 the experiment according to procedures approved by a local ethics committee. The data
1402 from 8 participants from Wenzhou site were excluded from analysis because their accuracy
1403 was close to chance (< 0.6). The results for the remaining 45 participants (33 female, age
1404 = 20.78 ± 1.76) were analyzed and reported.

1405 ***Design.***

1406 As in Experiment 3, a $2 \times 3 \times 2$ within-subject design was used. The first variable was
1407 self-relevance (self and stranger associations); the second variable was moral valence (good,
1408 neutral and bad associations); the third variable was the matching between shape and label
1409 (matching vs. non-match for the personal association). However, in this the task,
1410 participants only learn the association between two geometric shapes and two labels (self
1411 and other), i.e., only self-relevance were related to the task. The moral valence
1412 manipulation was achieved by embedding the personal label of the labels in the geometric
1413 shapes, see below. For simplicity, the trials where shapes where paired with self and with a
1414 word of “good person” inside were shorted as good-self condition, similarly, the trials where
1415 shapes paired with the self and with a word of “bad person” inside were shorted as bad-self
1416 condition. Hence, we also have six conditions: good-self, neutral-self, bad-self, good-other,
1417 neutral-other, and bad-other.

1418 ***Stimuli.***

1419 2 shapes were included (circle, square) and each appeared above a central fixation
1420 cross with the personal label appearing below. However, the shapes were not empty but
1421 with a two-Chinese-character word in the middle, the word was one of three labels with
1422 different moral valence: “good person”, “bad person” and “neutral person”. Before the
1423 experiment, participants learned the self/other association, and were informed to only
1424 response to the association between shapes’ configure and the labels below the fixation, but
1425 ignore the words within shapes. Besides the behavioral experiments, participants from

₁₄₂₆ Tsinghua community also finished questionnaires as Experiments 3, and participants from
₁₄₂₇ Wenzhou community finished a series of questionnaire as the other experiment finished in
₁₄₂₈ Wenzhou.

₁₄₂₉ ***Procedure.***

₁₄₃₀ The procedure was similar to Experiment 1. There were 6 blocks of trial, each with
₁₄₃₁ 120 trials for 2017 data. Due to procedure error, the data collected in 2015 in Tsinghua
₁₄₃₂ community only have 60 trials for each block, i.e., 30 trials per condition.

₁₄₃₃ As in study 3a, before each task, the instruction showed the meaning of each label to
₁₄₃₄ participants. The self-matching task and other-matching task were randomized between
₁₄₃₅ participants. Each participant finished 6 blocks, each have 120 trials.

₁₄₃₆ ***Data Analysis.***

₁₄₃₇ Same as experiment 3a.

₁₄₃₈ **Results.**

₁₄₃₉ ***NHST.***

₁₄₄₀ Figure 28 shows d prime and reaction times of experiment 3a. Less than 5% correct
₁₄₄₁ trials with less than 200ms reaction times were excluded.

₁₄₄₂ d prime.

₁₄₄₃ There was no evidence for the main effect of valence, $F(1.59, 69.94) = 2.34$,
₁₄₄₄ $MSE = 0.48$, $p = .115$, $\hat{\eta}_G^2 = .010$, but we found a main effect of self-relevance,
₁₄₄₅ $F(1, 44) = 0.00$, $MSE = 0.08$, $p = .994$, $\hat{\eta}_G^2 = .000$, as well as the interaction,
₁₄₄₆ $F(1.96, 86.41) = 0.53$, $MSE = 0.10$, $p = .585$, $\hat{\eta}_G^2 = .001$.

₁₄₄₇ We then conducted separated ANOVA for self-referential and other-referential trials.
₁₄₄₈ The valence effect was shown for the self-referential conditions, $F(1.75, 76.86) = 3.08$,
₁₄₄₉ $MSE = 0.25$, $p = .058$, $\hat{\eta}_G^2 = .017$. Post-hoc test revealed that the Good-Self condition
₁₄₅₀ (2.15 ± 0.12) was with greater d prime than Neutral condition (1.83 ± 0.12 , $t(34) = 3.36$,

¹⁴⁵¹ $p = 0.0031$), and Bad-self condition (1.87 ± 0.12), $t(34) = 2.955$, $p = 0.01$. There was
¹⁴⁵² difference between neutral and bad condition, $t(34) = -0.039$, $p = 0.914$. However, no effect
¹⁴⁵³ of valence was found for the other-referential condition $F(1.63, 71.50) = 1.07$, $MSE = 0.33$,
¹⁴⁵⁴ $p = .336$, $\hat{\eta}_G^2 = .006$.

¹⁴⁵⁵ *Reaction time.*

¹⁴⁵⁶ We found interaction between Matchness and Valence ($F(1.87, 82.50) = 18.58$,
¹⁴⁵⁷ $MSE = 1,291.12$, $p < .001$, $\hat{\eta}_G^2 = .023$) and then analyzed the matched trials and
¹⁴⁵⁸ nonmatch trials separately, as in previous experiments.

¹⁴⁵⁹ For the match trials, we found that the interaction between identity and valence,
¹⁴⁶⁰ $F(1.86, 81.84) = 5.22$, $MSE = 308.30$, $p = .009$, $\hat{\eta}_G^2 = .003$, as well as the main effect of
¹⁴⁶¹ valence $F(1.80, 79.37) = 11.04$, $MSE = 2,937.54$, $p < .001$, $\hat{\eta}_G^2 = .059$, but not the effect of
¹⁴⁶² identity $F(1, 44) = 0.23$, $MSE = 263.26$, $p = .632$, $\hat{\eta}_G^2 = .000$. As for the d prime, we
¹⁴⁶³ separated analyzed the self-referential and other-referential trials. For the Self-referential
¹⁴⁶⁴ trials, we found the main effect of valence, $F(1.74, 76.48) = 13.69$, $MSE = 1,732.08$,
¹⁴⁶⁵ $p < .001$, $\hat{\eta}_G^2 = .079$; for the other-referential trials, the effect of valence is weaker,
¹⁴⁶⁶ $F(1.87, 82.44) = 7.09$, $MSE = 1,527.43$, $p = .002$, $\hat{\eta}_G^2 = .043$. We then focused on the self
¹⁴⁶⁷ conditions: the good-self condition (713 ± 12) is faster than neutral- (776 ± 11.8), $t(34) =$
¹⁴⁶⁸ -7.396 , $p < .0001$, and bad-self (772 ± 10.1) conditions, $t(34) = -5.66$, $p < .0001$. But
¹⁴⁶⁹ there is not difference between neutral- and bad-self conditions, $t(34) = 0.481$, $p = 0.881$.

¹⁴⁷⁰ For the nonmatch trials, we didn't found any strong effect: identity, $F(1, 44) = 1.96$,
¹⁴⁷¹ $MSE = 319.47$, $p = .169$, $\hat{\eta}_G^2 = .001$, valence $F(1.69, 74.54) = 6.59$, $MSE = 886.19$,
¹⁴⁷² $p = .004$, $\hat{\eta}_G^2 = .010$, or interaction between the two $F(1.88, 82.57) = 0.31$, $MSE = 316.96$,
¹⁴⁷³ $p = .718$, $\hat{\eta}_G^2 = .000$.

¹⁴⁷⁴ **BGLM.**

¹⁴⁷⁵ *Signal detection theory analysis of accuracy.*

¹⁴⁷⁶ We found that the d prime is greater when shapes were associated with good self

1477 condition than with neutral self or bad self, but shapes associated with bad self and neutral
1478 self didn't show differences. comparing the self vs other under three condition revealed that
1479 shapes associated with good self is greater than with good other, but with a weak evidence.
1480 In contrast, for both neutral and bad valence condition, shapes associated with other had
1481 greater d prime than with self.

1482 *Reaction time.*

1483 In reaction times, we found that same trends in the match trials as in the RT: while
1484 the shapes associated with good self was greater than with good other (log mean diff =
1485 -0.02858, 95%HPD[-0.070898, 0.0154]), the direction is reversed for neutral and negative
1486 condition. see Figure 29

1487 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et
1488 al., 2020). We estimated separate drift rate (v), non-decision time (T_0), and boundary
1489 separation (a) for each condition. We found that the shapes tagged with good person has
1490 higher drift rate and higher boundary separation than shapes tagged with both neutral and
1491 bad person. Also, the shapes tagged with neutral person has a higher drift rate than
1492 shapes tagged with bad person, but not for the boundary separation. Finally, we found
1493 that shapes tagged with bad person had longer non-decision time (see figure 30)).

1494

Results

1495 **Effect of moral valence**

1496 In this part, we synthesized results from experiment 1a, 1b, 1c, 2, 5 and 6a. Data
1497 from 192 participants were included in these analyses. We found differences between
1498 positive and negative conditions on RT was Cohen's $d = -0.58 \pm 0.06$, 95% CI [-0.70 -0.47];
1499 on d' was Cohen's $d = 0.24 \pm 0.05$, 95% CI [0.15 0.34]. The effect was also observed
1500 between positive and neutral condition, RT: Cohen's $d = -0.44 \pm 0.10$, 95% CI [-0.63
1501 -0.25]; d' : Cohen's $d = 0.31 \pm 0.07$, 95% CI [0.16 0.45]. And the difference between neutral

1502 and bad conditions are not significant, RT: Cohen's $d = 0.15 \pm 0.07$, 95% CI [0.00 0.30];
1503 d' : Cohen's $d = 0.07 \pm 0.07$, 95% CI [-0.08 0.21]. See Figure 31 left panel.

1504 **Interaction between valence and self-reference**

1505 In this part, we combined the experiments that explicitly manipulated the
1506 self-reference and valence, which includes 3a, 3b, 6b, 7a, and 7b. For the positive versus
1507 negative contrast, data were from five experiments with 178 participants; for positive
1508 versus neutral and neutral versus negative contrasts, data were from three experiments (1509 3a, 3b, and 6b) with 108 participants.

1510 In most of these experiments, the interaction between self-reference and valence was
1511 significant (see results of each experiment in supplementary materials). In the
1512 mini-meta-analysis, we analyzed the valence effect for self-referential condition and
1513 other-referential condition separately.

1514 For the self-referential condition, we found the same pattern as in the first part of
1515 results. That is we found significant differences between positive and neutral as well as
1516 positive and negative, but not neutral and negative. The effect size of RT between positive
1517 and negative is Cohen's $d = -0.89 \pm 0.12$, 95% CI [-1.11 -0.66]; on d' was Cohen's $d = 0.61$
1518 ± 0.09 , 95% CI [0.44 0.78]. The effect was also observed between positive and neutral
1519 condition, RT: Cohen's $d = -0.76 \pm 0.13$, 95% CI [-1.01 -0.50]; d' : Cohen's $d = 0.69 \pm$
1520 0.14, 95% CI [0.42 0.96]. And the difference between neutral and bad conditions are not
1521 significant, RT: Cohen's $d = 0.03 \pm 0.13$, 95% CI [-0.22 0.29]; d' : Cohen's $d = 0.08 \pm 0.08$,
1522 95% CI [-0.07 0.24]. See Figure 31 the middle panel.

1523 For the other-referential condition, we found that only the difference between positive
1524 and negative on RT was significant, all the other conditions were not. The effect size of RT
1525 between positive and negative is Cohen's $d = -0.28 \pm 0.05$, 95% CI [-0.38 -0.17]; on d' was
1526 Cohen's $d = -0.02 \pm 0.08$, 95% CI [-0.17 0.13]. The effect was not observed between

1527 positive and neutral condition, RT: Cohen's $d = -0.12 \pm 0.10$, 95% CI [-0.31 0.06]; d' :
1528 Cohen's $d = 0.01 \pm 0.08$, 95% CI [-0.16 0.17]. And the difference between neutral and bad
1529 conditions are not significant, RT: Cohen's $d = 0.14 \pm 0.09$, 95% CI [-0.03 0.31]; d' :
1530 Cohen's $d = 0.05 \pm 0.07$, 95% CI [-0.08 0.18]. See Figure 31 right panel.

1531 **Generalizability of the valence effect**

1532 In this part, we reported the results from experiment 4 in which either moral valence
1533 or self-reference were manipulated as task-irrelevant stimuli.

1534 For experiment 4a, when self-reference was the target and moral valence was
1535 task-irrelevant, we found that only under the implicit self-referential condition, i.e., when
1536 the moral words were presented as task irrelevant stimuli, there was the main effect of
1537 valence and interaction between valence and reference for both d prime and RT (See
1538 supplementary results for the detailed statistics). For d prime, we found good-self
1539 condition (2.55 ± 0.86) had higher d prime than bad-self condition (2.38 ± 0.80); good self
1540 condition was also higher than neutral self (2.45 ± 0.78) but there was not statistically
1541 significant, while the neutral-self condition was higher than bad self condition and not
1542 significant neither. For reaction times, good-self condition (654.26 ± 67.09) were faster
1543 relative to bad-self condition (665.64 ± 64.59), and over neutral-self condition ($664.26 \pm$
1544 64.71). The difference between neutral-self and bad-self conditions were not significant.
1545 However, for the other-referential condition, there was no significant differences between
1546 different valence conditions. See Figure 32.

1547 For experiment 4b, when valence was the target and the identity was task-irrelevant,
1548 we found a strong valence effect (see supplementary results and Figure 33, Figure 34).

1549 In this experiment, the advantage of good-self condition can only be disentangled by
1550 comparing the self-referential and other-referential conditions. Therefore, we calculated the
1551 differences between the valence effect under self-referential and other referential conditions

1552 and used the weighted variance as the variance of this differences. We found this
1553 modulation effect on RT. The valence effect of RT was stronger in self-referential than
1554 other-referential for the Good vs. Neutral condition (-0.33 ± 0.01), and to a less extent the
1555 Good vs. Bad condition (-0.17 ± 0.01). While the size of the other effect's CI included
1556 zero, suggesting those effects didn't differ from zero. See Figure 35.

1557 **Specificity of valence effect**

1558 In this part, we analyzed the results from experiment 5, which included positive,
1559 neutral, and negative valence from four different domains: morality, emotion, aesthetics of
1560 human, and aesthetics of scene. We found interaction between valence and domain for both
1561 d prime and RT (match trials). A common pattern appeared in all four domains: each
1562 domain showed a binary results instead of gradient on both d prime and RT. For morality,
1563 aesthetics of human, and aesthetics of scene, there was a positivity effect where the positive
1564 conditions had advantages over both neutral (greater d prime and faster RT), while neutral
1565 and negative conditions didn't differ from each other. But for the emotional stimuli, there
1566 was a reversed negativity effect: positive and neutral conditions were not significantly
1567 different from each other but both had advantage over negative conditions. See
1568 supplementary materials for detailed statistics. Also note that the effect size in moral
1569 domain is smaller than the aesthetic domains (beauty of people and beauty of scene). See
1570 Figure 36.

1571 **Self-reported personal distance**

1572 See Figure 37.

1573 **Correlation analyses**

1574 The reliability of questionnaires can be found in (Liu et al., 2020). We calculated the
1575 correlation between the data from behavioral task and the questionnaire data. First, we
1576 calculated the score for each scale based on their structure and factor loading, instead of
1577 sum score (McNeish & Wolf, 2020). Then, we used SEM to estimate the correlation
1578 because it can include measurement model and statistical model in a unified framework.

1579 To make sure that what we found were not false positive, we used two method to
1580 ensure the robustness of our analysis. first, we split the data into two half: the data with
1581 self and without, then, we used the conditional random forest to find the robust correlation
1582 in the exploratory data (with self reference) that can be replicated in the confirmatory data
1583 (without the self reference). The robust correlation were then analyzed using SEM

1584 Instead of use the exploratory correlation analysis, we used a more principled way to
1585 explore the correlation between parameter of HDDM (v , t , and a) and scale scores and
1586 person distance.

1587 We didn't find the correlation between scale scores and the parameters of HDDM,
1588 but found weak correlation between personal distance and the parameter estimated from
1589 Good and neutral conditions.

1590 First, boundary separation (a) of moral good condition was correlated with both
1591 Self-Bad distance ($r = 0.198$, 95% CI [], $p = 0.0063$) and Neutral-Bad distance
1592 ($r = 0.1571$, 95% CI [], $p = 0.031$). At the same time, the non-decision time is negatively
1593 correlated with Self-Bad distance ($r = 0.169$, 95% CI [], $p = 0.0197$). See Figure 38.

1594 Second, we found the boundary separation of neutral condition is positively
1595 correlated with the personal distance between self and good distance ($r = 0.189$, 95% CI [],
1596 $p = 0.036$), but negatively correlated with self-neutral distance($r = -0.183$, 95% CI [],
1597 $p = 0.042$). Also, the drift rate of the neutral condition is positively correlated with the
1598 Self-Bad distance ($r = 0.177$, 95% CI [], $p = 0.048$).a. See figure 39

1599 We also explored the correlation between behavioral data and questionnaire scores
1600 separately for experiments with and without self-referential, however, the sample size is
1601 very low for some conditions.

1602 **Discussion**

1603 **References**

- 1604 Abele, A. E., Ellemers, N., Fiske, S. T., Koch, A., & Yzerbyt, V. (2020). Navigating the
1605 social world: Toward an integrated framework for evaluating self, individuals, and
1606 groups. *Psychological Review*. <https://doi.org/10.1037/rev0000262>
- 1607 Amodio, D. M. (2019). Social cognition 2.0: An interactive memory systems account.
1608 *Trends in Cognitive Sciences*, 23(1), 21–33.
1609 <https://doi.org/10.1016/j.tics.2018.10.002>
- 1610 Anderson, E., Siegel, E. H., Bliss-Moreau, E., & Barrett, L. F. (2011). The visual impact
1611 of gossip. *Science*, 332(6036), 1446–1448. <https://doi.org/10.1126/science.1201574>
- 1612 Brainard, D. H. (1997). The psychophysics toolbox. *Spatial Vision*, 10(4), 433–436.
1613 Journal Article.
- 1614 Bürkner, P.-C. (2017). Brms: An r package for bayesian multilevel models using stan.
1615 *Journal of Statistical Software; Vol 1, Issue 1 (2017)*. Journal Article. Retrieved
1616 from
1617 <https://www.jstatsoft.org/v080/i01%0Ahttp://dx.doi.org/10.18637/jss.v080.i01>
- 1618 Carlson, R. W., Maréchal, M. A., Oud, B., Fehr, E., & Crockett, M. J. (2020). Motivated
1619 misremembering of selfish decisions. *Nature Communications*, 11(1), 2100.
1620 <https://doi.org/10.1038/s41467-020-15602-4>
- 1621 Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M., ...
1622 Riddell, A. (2017). Stan: A probabilistic programming language. *Journal of*

- 1623 *Statistical Software*, 76(1). Journal Article. <https://doi.org/10.18637/jss.v076.i01>
- 1624 Cooper, H., Hedges, L. V., & Valentine, J. C. (2009). *The handbook of research synthesis
1625 and meta-analysis* (2nd ed.). Book, New York: Sage.
- 1626 DeCarlo, L. T. (1998). Signal detection theory and generalized linear models. *Psychological
1627 Methods*, 3(2), 186–205. Journal Article. <https://doi.org/10.1037/1082-989X.3.2.186>
- 1628 Dunlea, J. P., & Heiphetz, L. (2020). Children's and adults' understanding of punishment
1629 and the criminal justice system. *Journal of Experimental Social Psychology*, 87,
1630 103913. <https://doi.org/10.1016/j.jesp.2019.103913>
- 1631 Dzhelyova, M., Perrett, D. I., & Jentzsch, I. (2012). Temporal dynamics of trustworthiness
1632 perception. *Brain Research*, 1435, 81–90.
1633 <https://doi.org/10.1016/j.brainres.2011.11.043>
- 1634 Enock, F., Sui, J., Hewstone, M., & Humphreys, G. W. (2018). Self and team prioritisation
1635 effects in perceptual matching: Evidence for a shared representation. *Acta
1636 Psychologica*, 182, 107–118. <https://doi.org/10.1016/j.actpsy.2017.11.011>
- 1637 Faul, F., Erdfelder, E., Buchner, A., & Lang, A.-G. (2009). Statistical power analyses using
1638 g*Power 3.1: Tests for correlation and regression analyses. *Behavior Research
1639 Methods*, 41(4), 1149–1160. <https://doi.org/10.3758/BRM.41.4.1149>
- 1640 Firestone, C., & Scholl, B. J. (2015). Enhanced visual awareness for morality and pajamas?
1641 Perception vs. Memory in “top-down” effects. *Cognition*, 136, 409–416.
1642 <https://doi.org/10.1016/j.cognition.2014.10.014>
- 1643 Freeman, J. B., & Ambady, N. (2011). A dynamic interactive theory of person construal.
1644 *Psychological Review*, 118(2), 247–279. <https://doi.org/10.1037/a0022327>
- 1645 Gantman, A. P., & Van Bavel, J. J. (2014). The moral pop-out effect: Enhanced
1646 perceptual awareness of morally relevant stimuli. *Cognition*, 132(1), 22–29.
1647 <https://doi.org/10.1016/j.cognition.2014.02.007>

- 1648 Goh, J. X., Hall, J. A., & Rosenthal, R. (2016). Mini meta-analysis of your own studies:
1649 Some arguments on why and a primer on how. *Social and Personality Psychology
1650 Compass*, 10(10), 535–549. Journal Article. <https://doi.org/10.1111/spc3.12267>
- 1651 Goodwin, G. P. (2015). Moral character in person perception. *Current Directions in
1652 Psychological Science*, 24(1), 38–44. <https://doi.org/10.1177/0963721414550709>
- 1653 Goodwin, G. P., Piazza, J., & Rozin, P. (2014). Moral character predominates in person
1654 perception and evaluation. *Journal of Personality and Social Psychology*, 106(1),
1655 148–168. <https://doi.org/10.1037/a0034726>
- 1656 Henrich, J., Heine, S. J., & Norenzayan, A. (2010). The weirdest people in the world?
1657 *Behavioral and Brain Sciences*, 33(2), 61–83.
1658 <https://doi.org/10.1017/S0140525X0999152X>
- 1659 Hofmann, W., Wisneski, D. C., Brandt, M. J., & Skitka, L. J. (2014). Morality in everyday
1660 life. *Science*, 345(6202), 1340–1343. <https://doi.org/10.1126/science.1251560>
- 1661 Hu, C.-P., Lan, Y., Macrae, C. N., & Sui, J. (2020). Good me bad me: Does valence
1662 influence self-prioritization during perceptual decision-making? *Collabra:
1663 Psychology*, 6(1), 20. Journal Article. <https://doi.org/10.1525/collabra.301>
- 1664 Juechems, K., & Summerfield, C. (2019). Where does value come from? *Trends in
1665 Cognitive Sciences*, 23(10), 836–850. <https://doi.org/10.1016/j.tics.2019.07.012>
- 1666 Kirby, K. N., & Gerlanc, D. (2013). BootES: An r package for bootstrap confidence
1667 intervals on effect sizes. *Behavior Research Methods*, 45(4), 905–927.
1668 <https://doi.org/10.3758/s13428-013-0330-5>
- 1669 Klein, N., & Epley, N. (2016). Maybe holier, but definitely less evil, than you: Bounded
1670 self-righteousness in social judgment. *Journal of Personality and Social Psychology*,
1671 110(5), 660–674. <https://doi.org/10.1037/pspa0000050>
- 1672 Liu, Q., Wang, F., Yan, W., Peng, K., Sui, J., & Hu, C.-P. (2020). Questionnaire data from

- 1673 the revision of a chinese version of free will and determinism plus scale. *Journal of*
1674 *Open Psychology Data*, 8(1), 1. Journal Article. <https://doi.org/10.5334/jopd.49/>
- 1675 Matzke, D., & Wagenmakers, E.-J. (2009). Psychological interpretation of the ex-gaussian
1676 and shifted wald parameters: A diffusion model analysis. *Psychonomic Bulletin &*
1677 *Review*, 16(5), 798–817. <https://doi.org/10.3758/PBR.16.5.798>
- 1678 McHugh, C., McGann, M., Igou, E. R., & Kinsella, E. (2019). *Moral judgment as*
1679 *categorization (MJAC)*. PsyArXiv. <https://doi.org/10.31234/osf.io/72dzp>
- 1680 McNeish, D., & Wolf, M. G. (2020). Thinking twice about sum scores. *Behavior Research*
1681 *Methods*. <https://doi.org/10.3758/s13428-020-01398-0>
- 1682 Monin, B., & Jordan, A. H. (2009). The dynamic moral self: A social psychological
1683 perspective. In *Personality, identity, and character: Explorations in moral*
1684 *psychology* (pp. 341–354). New York, NY, US: Cambridge University Press.
1685 <https://doi.org/10.1017/CBO9780511627125.016>
- 1686 Pelli, D. G. (1997). The videotoolbox software for visual psychophysics: Transforming
1687 numbers into movies. *Spatial Vision*, 10(4), 437–442. Journal Article.
- 1688 Reicher, S., & Hopkins, N. (2016). Perception, action, and the social dynamics of the
1689 variable self. *Psychological Inquiry*, 27(4), 341–347.
1690 <https://doi.org/10.1080/1047840X.2016.1217584>
- 1691 Rouder, J. N., & Lu, J. (2005). An introduction to bayesian hierarchical models with an
1692 application in the theory of signal detection. *Psychonomic Bulletin & Review*,
1693 12(4), 573–604. Journal Article. <https://doi.org/10.3758/bf03196750>
- 1694 Rousselet, G. A., & Wilcox, R. R. (2019). Reaction times and other skewed distributions:
1695 Problems with the mean and the median. *Meta-Psychology*. preprint.
1696 <https://doi.org/10.1101/383935>
- 1697 Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2013). Life after p-hacking. Conference

- 1698 Proceedings. <https://doi.org/10.2139/ssrn.2205186>
- 1699 Stanislaw, H., & Todorov, N. (1999). Calculation of signal detection theory measures.
1700 *Behavior Research Methods, Instruments, & Computers*, 31(1), 137–149. Journal
1701 Article. <https://doi.org/10.3758/BF03207704>
- 1702 Stanley, M. L., & De Brigard, F. (2019). Moral memories and the belief in the good self.
1703 *Current Directions in Psychological Science*, 28(4), 387–391.
1704 <https://doi.org/10.1177/0963721419847990>
- 1705 Stein, T., Grubb, C., Bertrand, M., Suh, S. M., & Verosky, S. C. (2017). No impact of
1706 affective person knowledge on visual awareness: Evidence from binocular rivalry and
1707 continuous flash suppression. *Emotion*, 17(8), 1199–1207.
1708 <https://doi.org/10.1037/emo0000305>
- 1709 Stolier, R. M., & Freeman, J. B. (2016). Functional and temporal considerations for
1710 top-down influences in social perception. *Psychological Inquiry*, 27(4), 352–357.
1711 <https://doi.org/10.1080/1047840X.2016.1216034>
- 1712 Strohminger, N., Knobe, J., & Newman, G. (2017). The true self: A psychological concept
1713 distinct from the self: *Perspectives on Psychological Science*.
1714 <https://doi.org/10.1177/1745691616689495>
- 1715 Sui, J., He, X., & Humphreys, G. W. (2012). Perceptual effects of social salience: Evidence
1716 from self-prioritization effects on perceptual matching. *Journal of Experimental
1717 Psychology: Human Perception and Performance*, 38(5), 1105–1117. Journal
1718 Article. <https://doi.org/10.1037/a0029792>
- 1719 Tappin, B. M., & McKay, R. T. (2017). The illusion of moral superiority. *Social
1720 Psychological and Personality Science*, 8(6), 623–631.
1721 <https://doi.org/10.1177/1948550616673878>
- 1722 Turner, J. C., Oakes, P. J., Haslam, S. A., & McGarty, C. (1994). Self and collective:

- 1723 Cognition and social context. *Personality and Social Psychology Bulletin*, 20(5),
1724 454–463. <https://doi.org/10.1177/0146167294205002>
- 1725 Uhlmann, E. L., Pizarro, D. A., & Diermeier, D. (2015). A person-centered approach to
1726 moral judgment: *Perspectives on Psychological Science*.
1727 <https://doi.org/10.1177/1745691614556679>
- 1728 Verosky, S. C., & Todorov, A. (2010). Differential neural responses to faces physically
1729 similar to the self as a function of their valence. *NeuroImage*, 49(2), 1690–1698.
1730 <https://doi.org/10.1016/j.neuroimage.2009.10.017>
- 1731 Wiecki, T. V., Sofer, I., & Frank, M. J. (2013). HDDM: Hierarchical bayesian estimation of
1732 the drift-diffusion model in python. *Frontiers in Neuroinformatics*, 7.
1733 <https://doi.org/10.3389/fninf.2013.00014>
- 1734 Willis, J., & Todorov, A. (2006). First impressions: Making up your mind after a 100-ms
1735 exposure to a face. *Psychological Science*, 17(7), 592–598.
1736 <https://doi.org/10.1111/j.1467-9280.2006.01750.x>
- 1737 Xiao, Y. J., Coppin, G., & Bavel, J. J. V. (2016). Perceiving the world through
1738 group-colored glasses: A perceptual model of intergroup relations. *Psychological
1739 Inquiry*, 27(4), 255–274. <https://doi.org/10.1080/1047840X.2016.1199221>

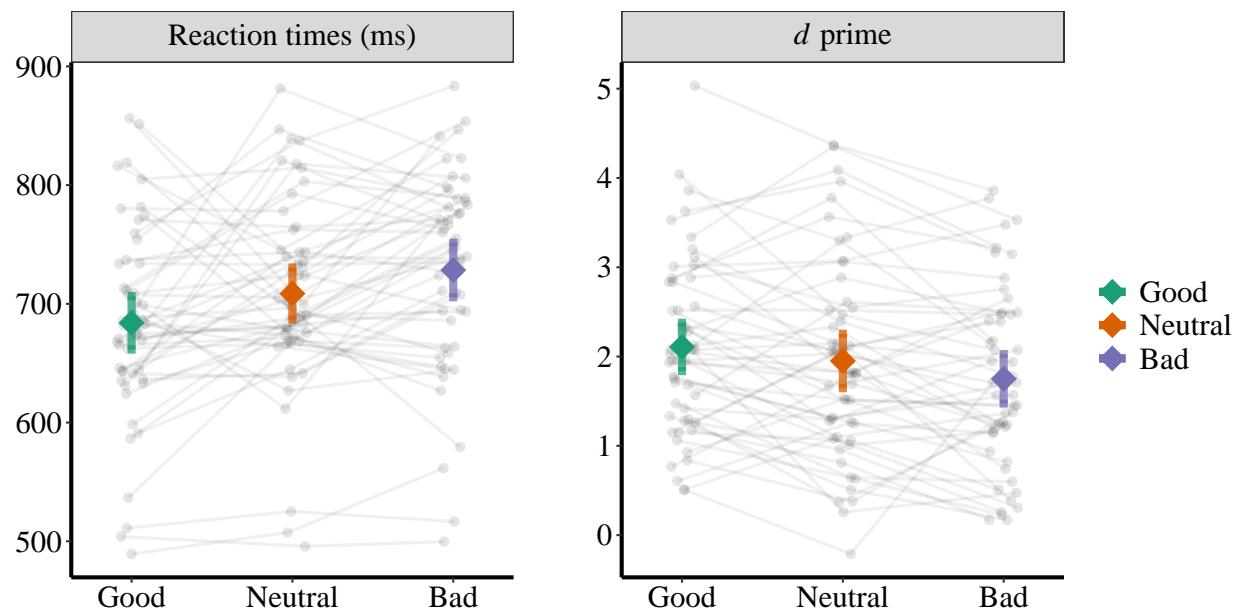


Figure 1. RT and d prime of Experiment 1a.

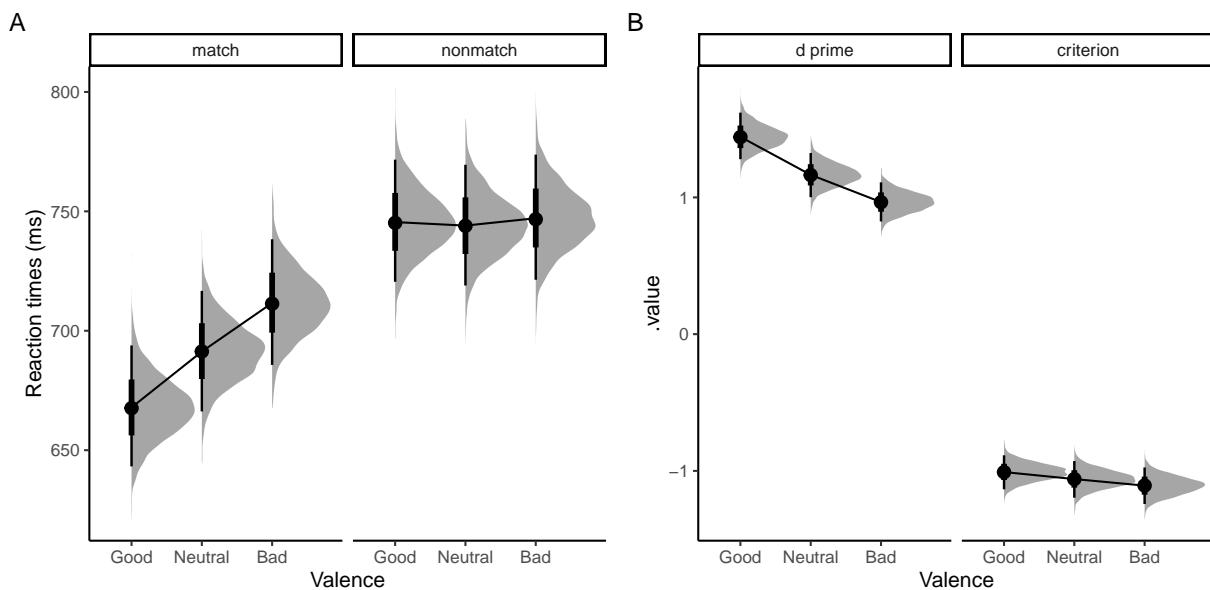


Figure 2. Exp1a: Results of Bayesian GLM analysis.

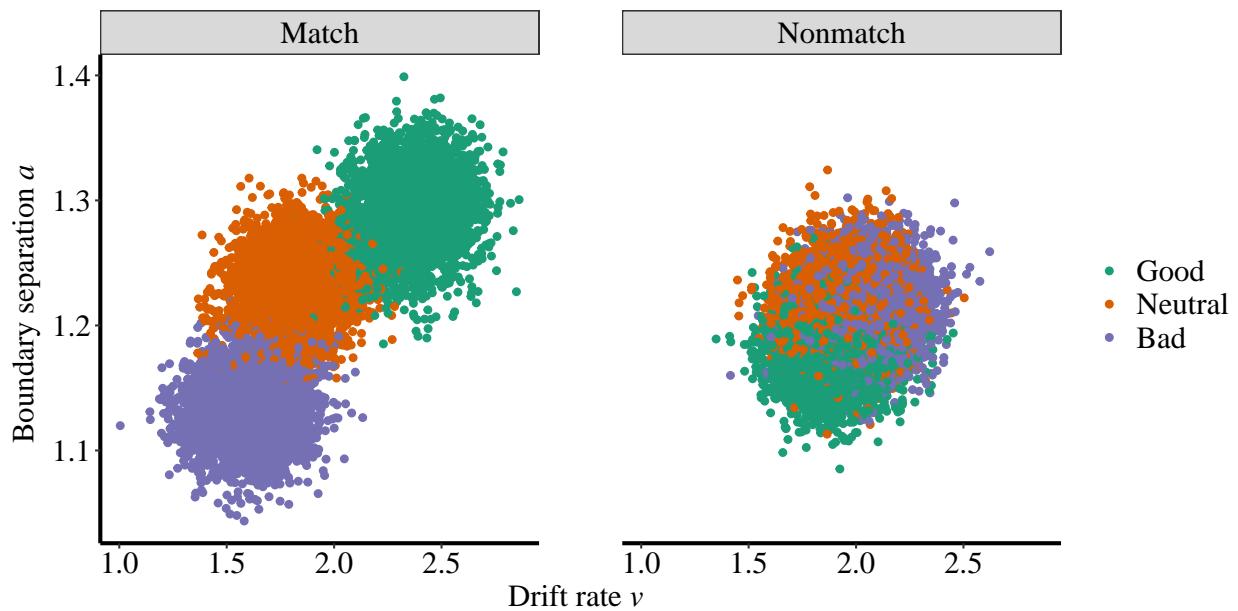


Figure 3. Exp1a: Results of HDDM.

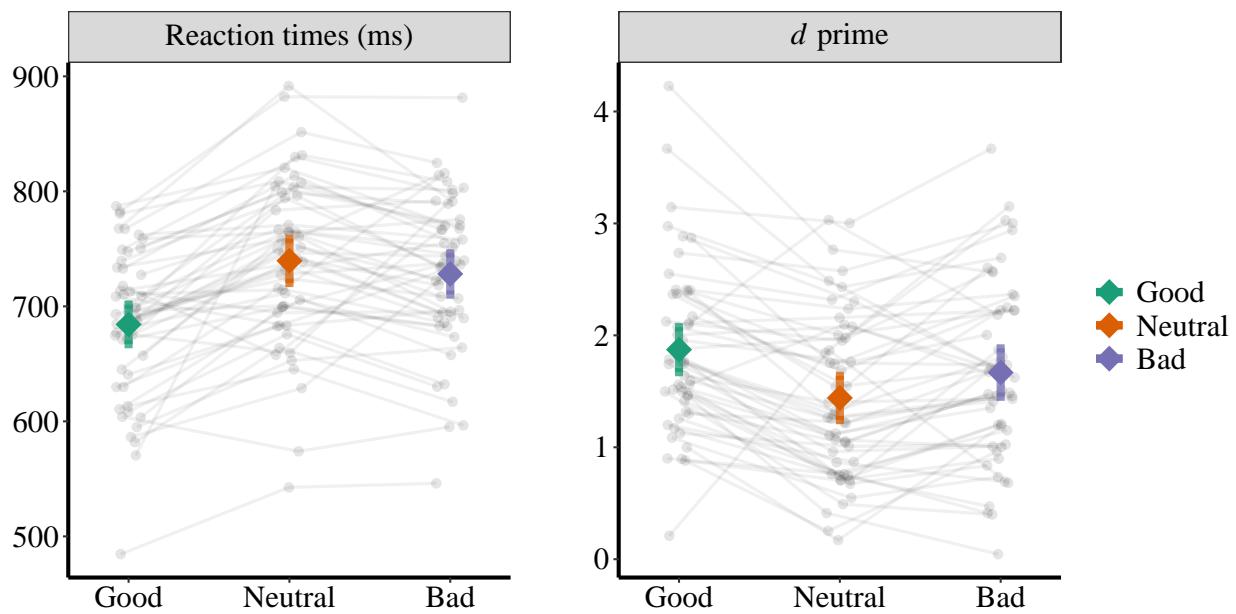


Figure 4. RT and d' of Experiment 1b.

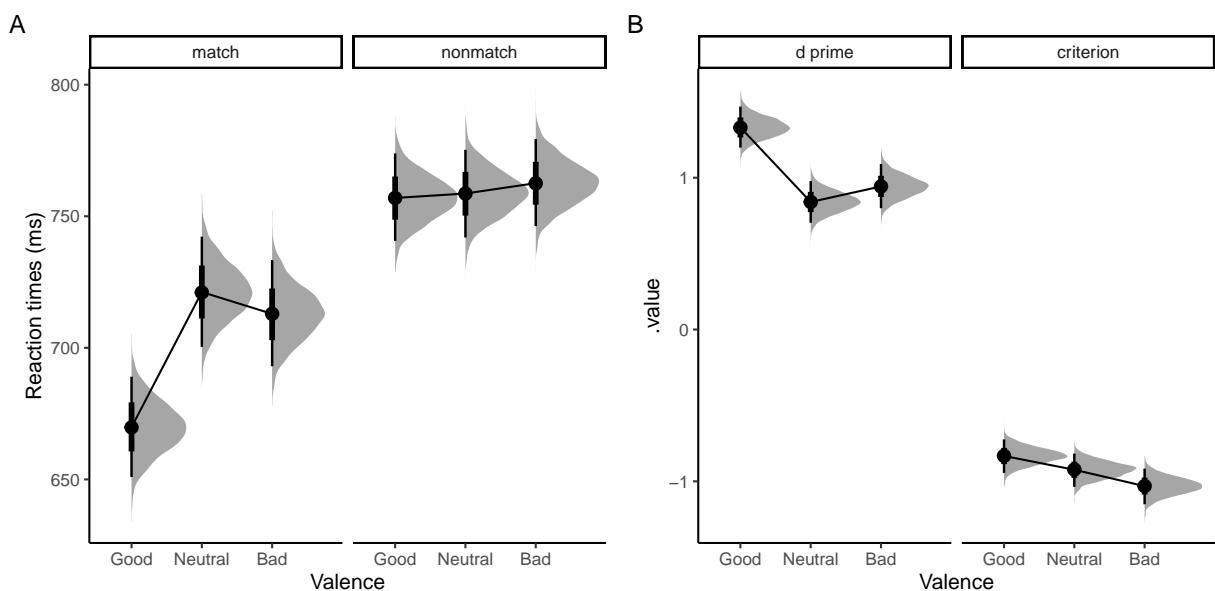


Figure 5. Exp1b: Results of Bayesian GLM analysis.

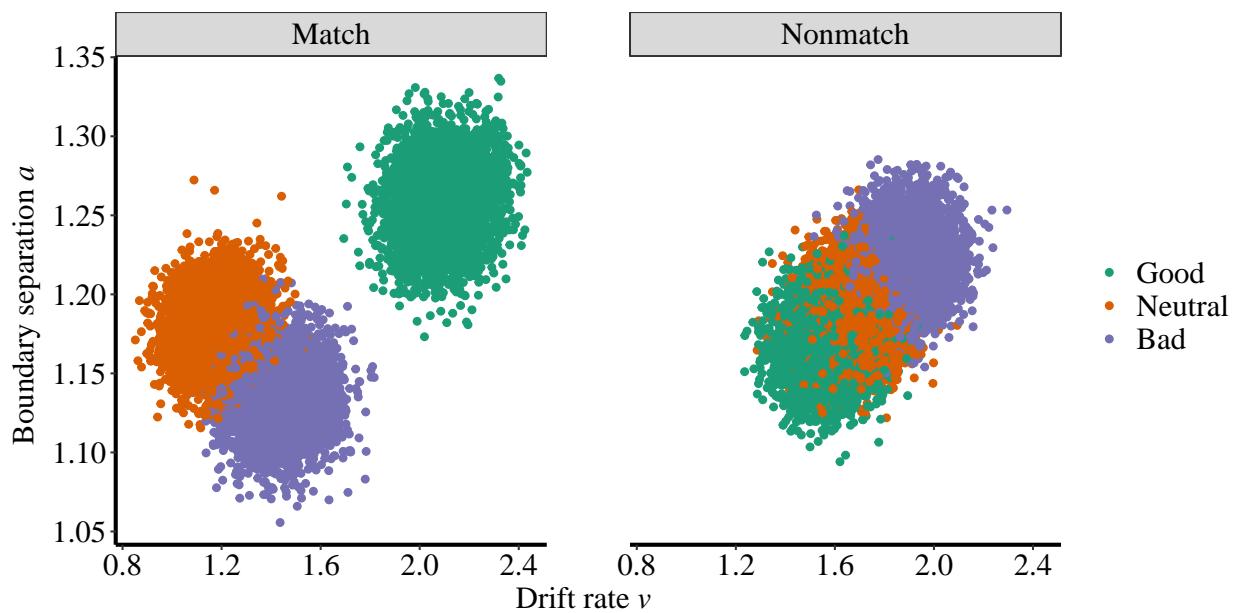


Figure 6. Exp1b: Results of HDDM.

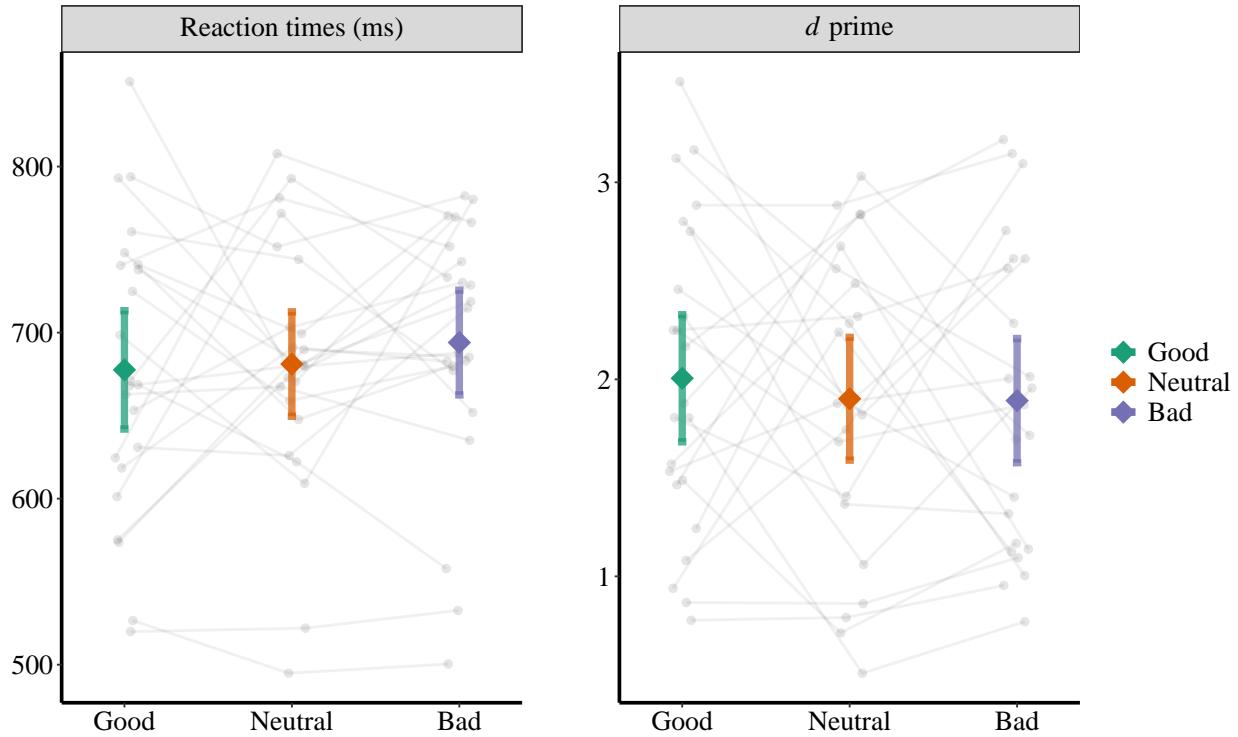


Figure 7. RT and d' prime of Experiment 1c.

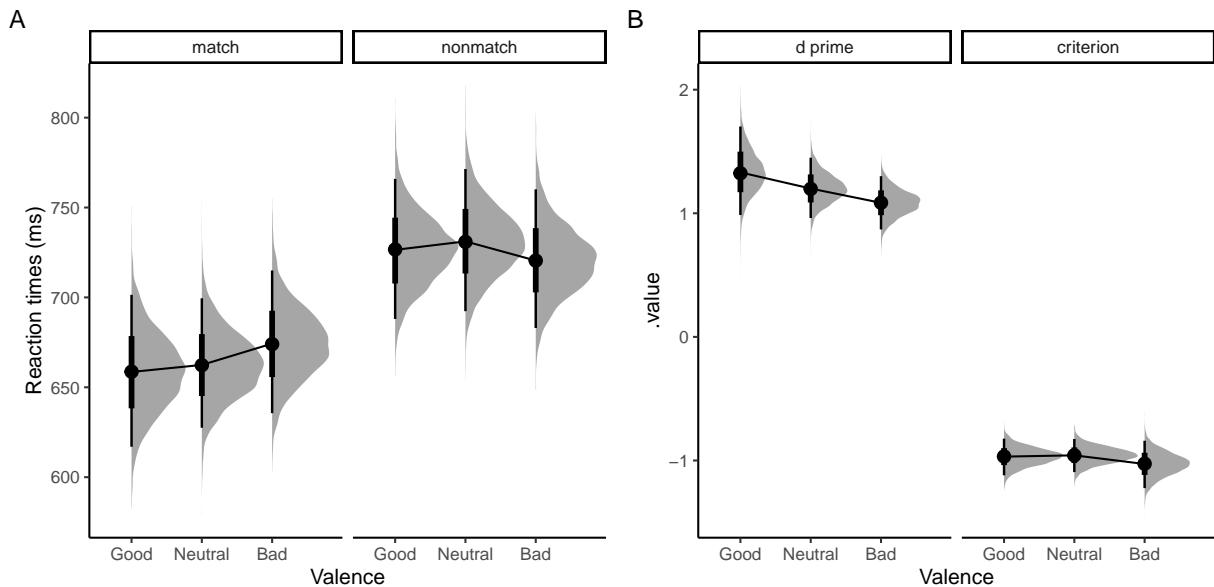


Figure 8. Exp1c: Results of Bayesian GLM analysis.

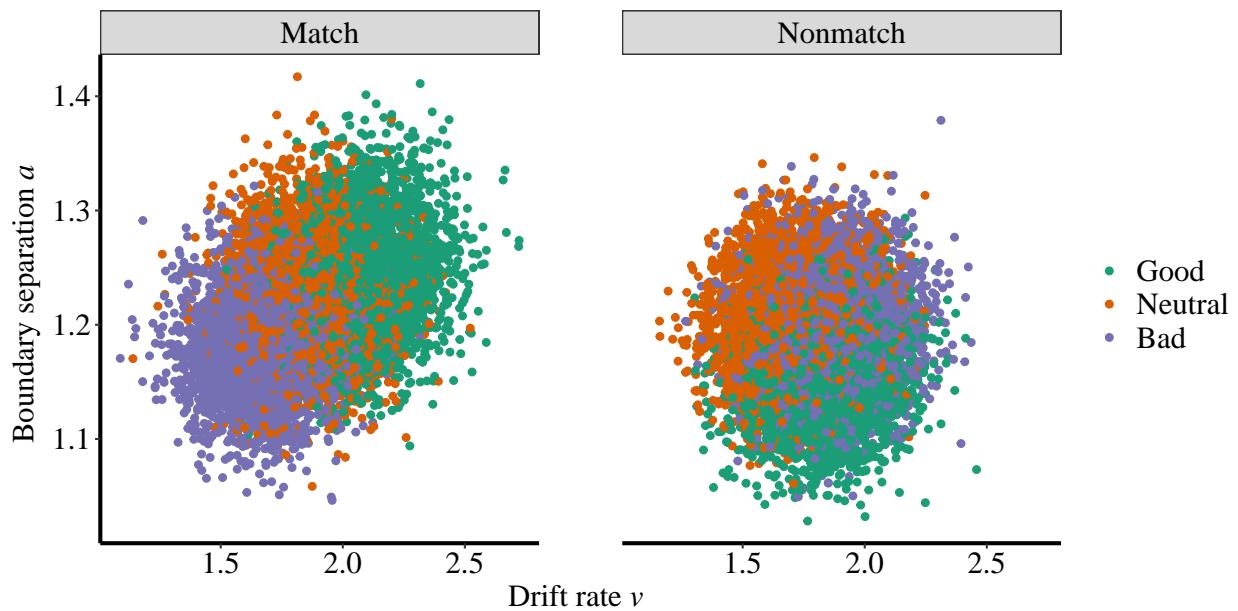


Figure 9. Exp1c: Results of HDDM.

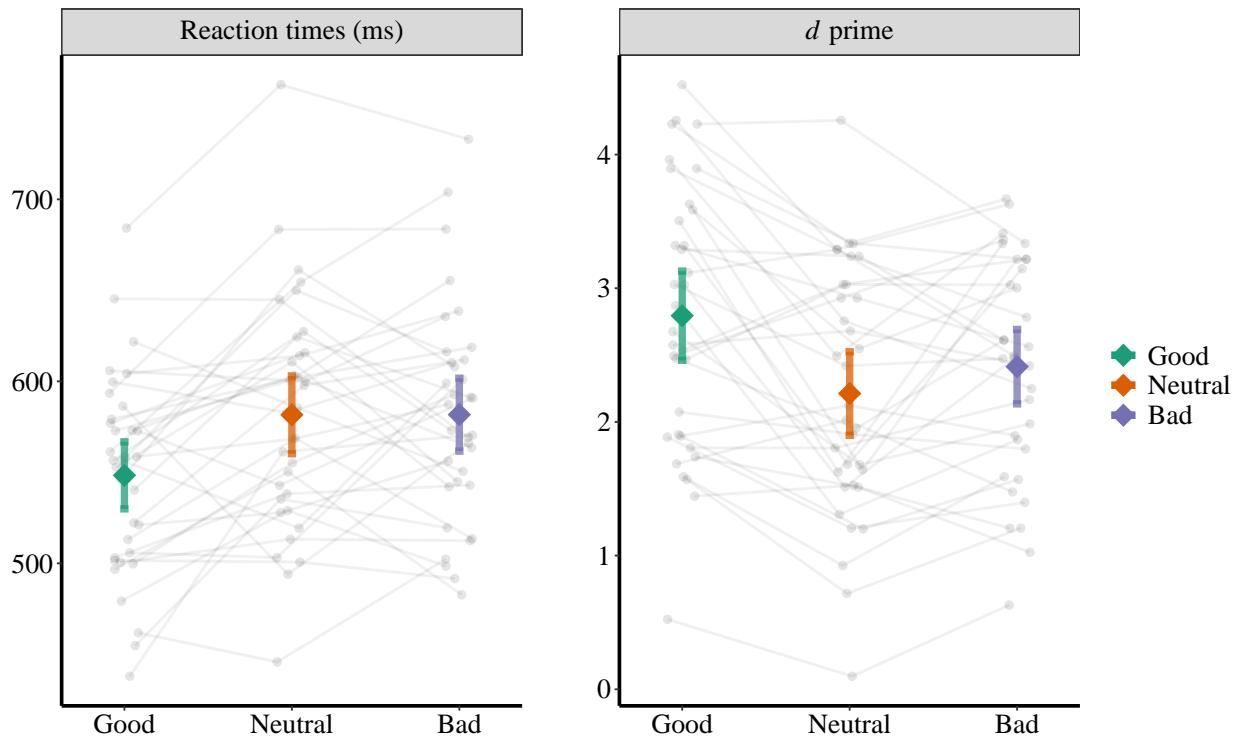


Figure 10. RT and d' of Experiment 2.

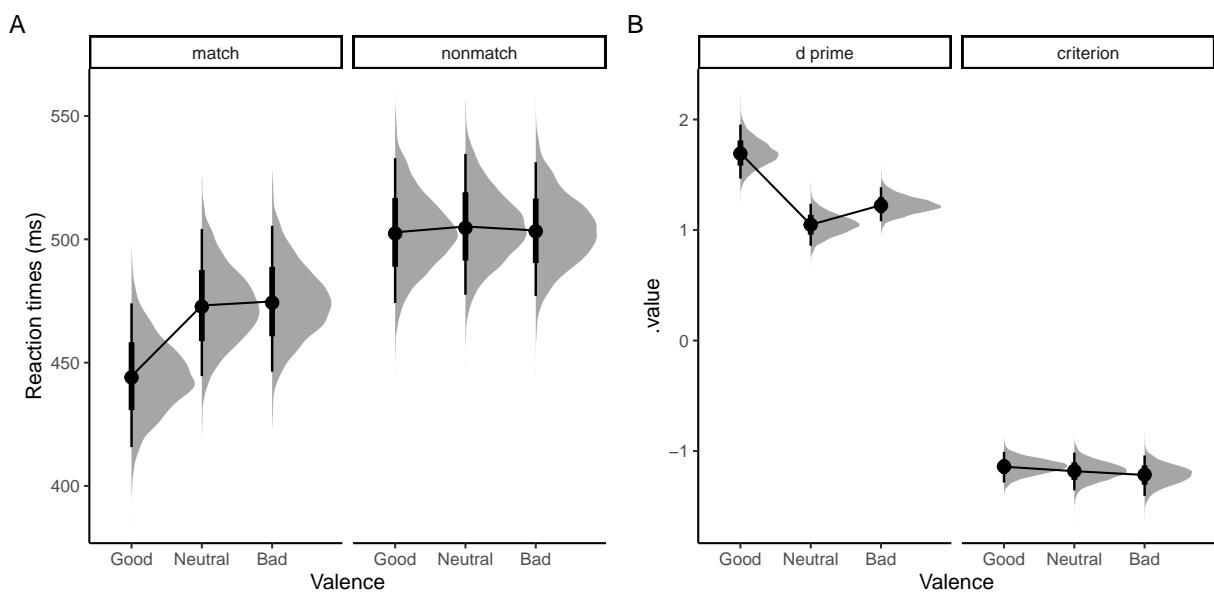


Figure 11. Exp2: Results of Bayesian GLM analysis.

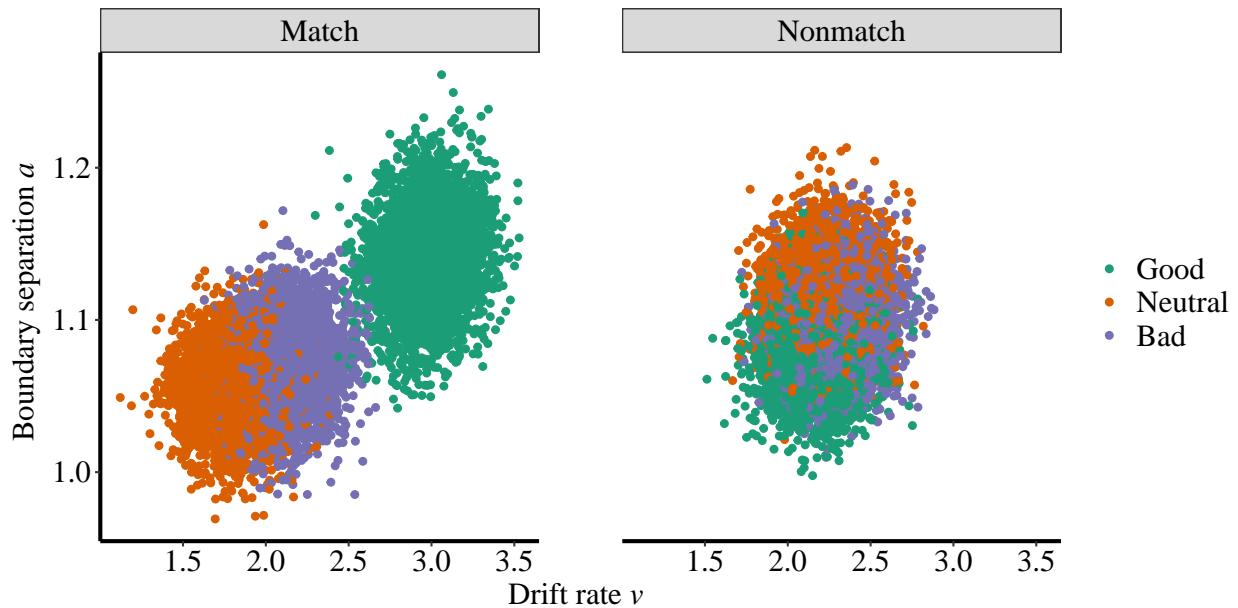


Figure 12. Exp2: Results of HDDM.

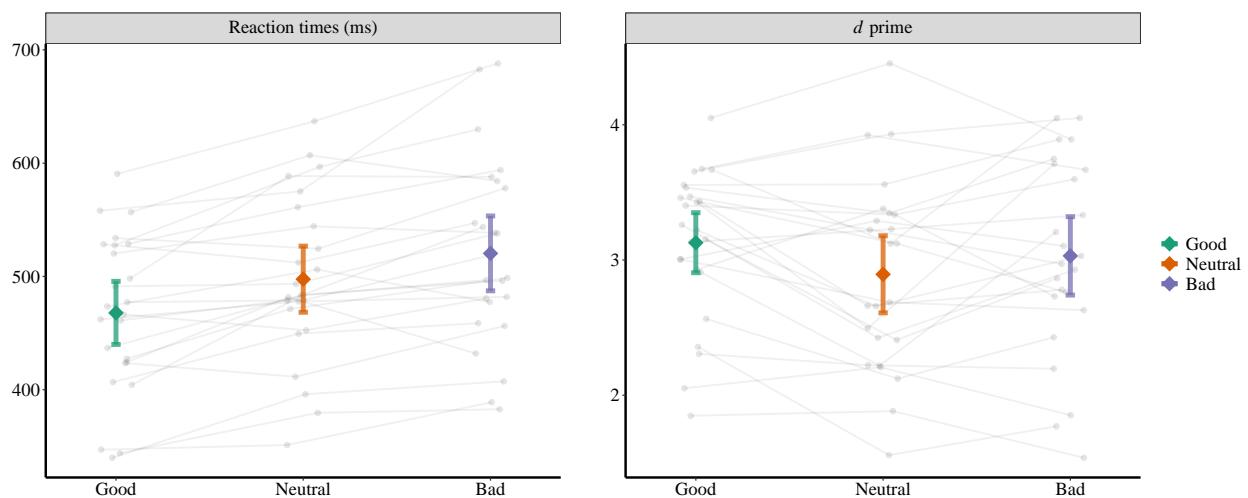


Figure 13. RT and d' of Experiment 6a.

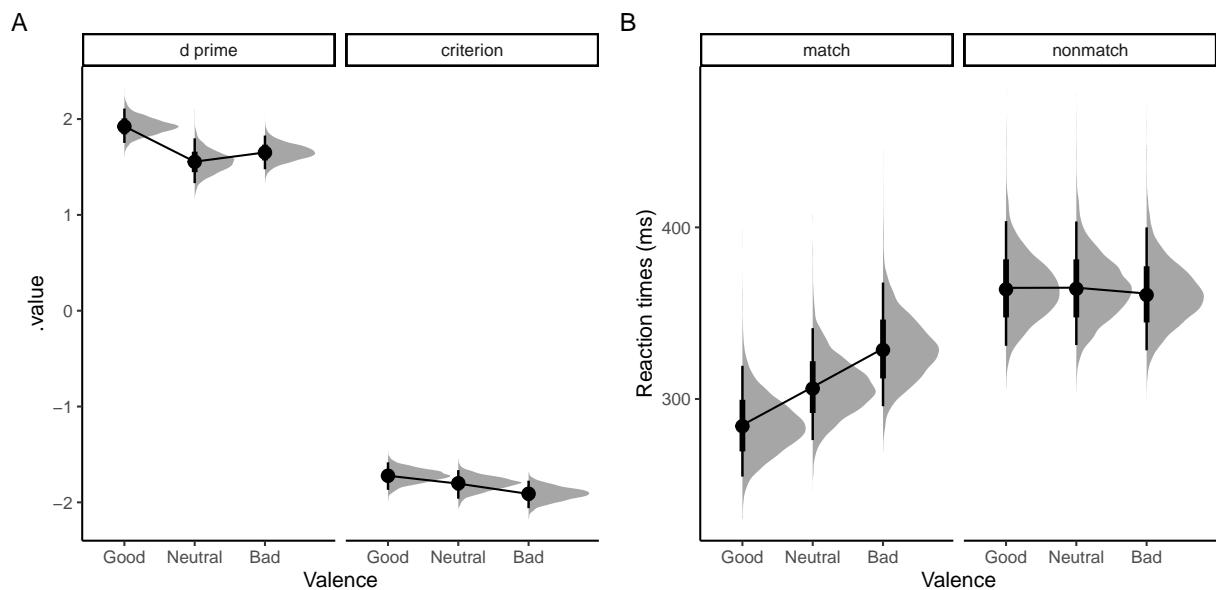


Figure 14. Exp6a: Results of Bayesian GLM analysis.

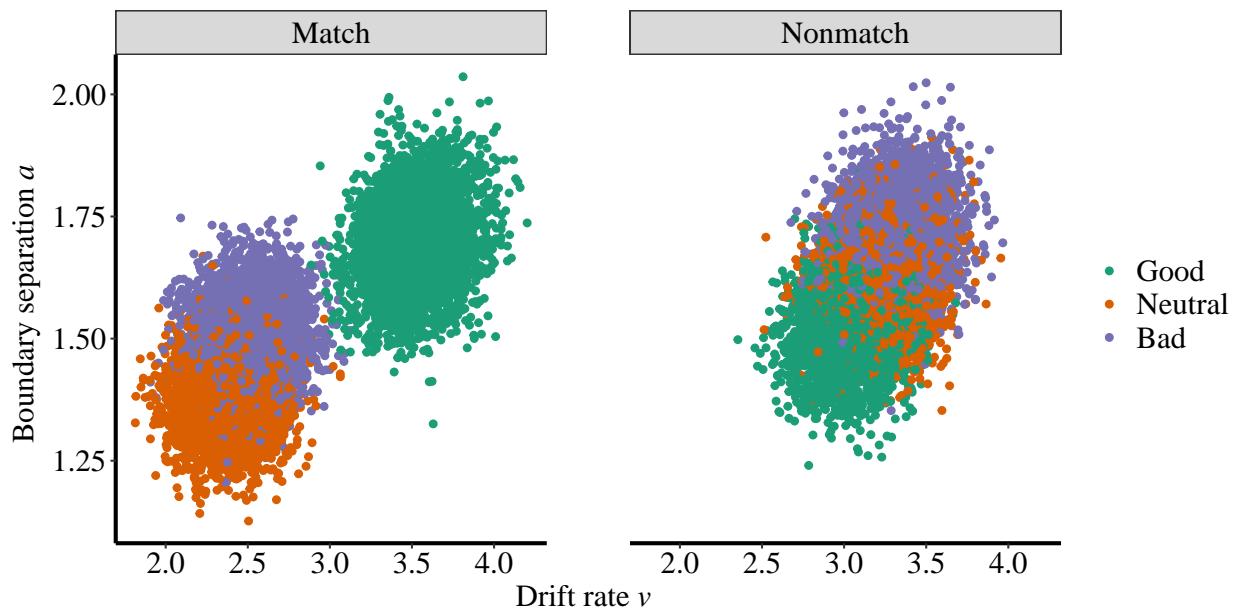


Figure 15. exp6a: Results of HDDM.

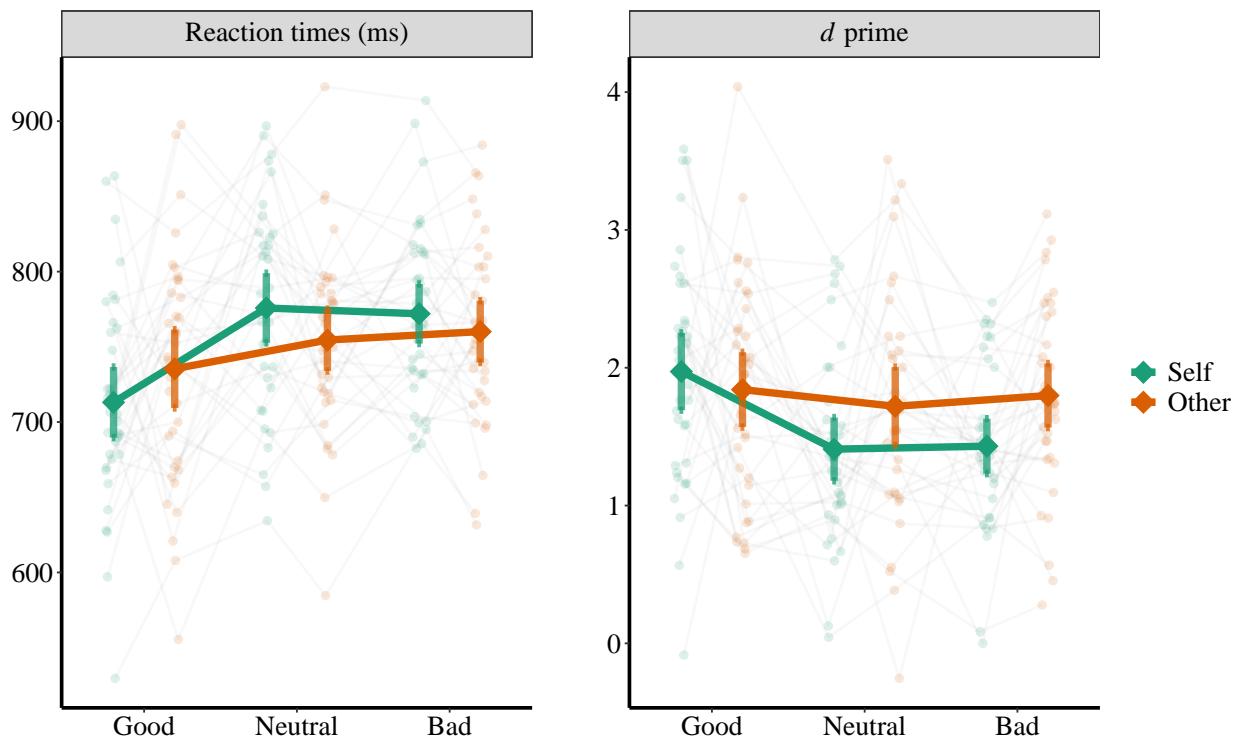


Figure 16. RT and d' prime of Experiment 3a.

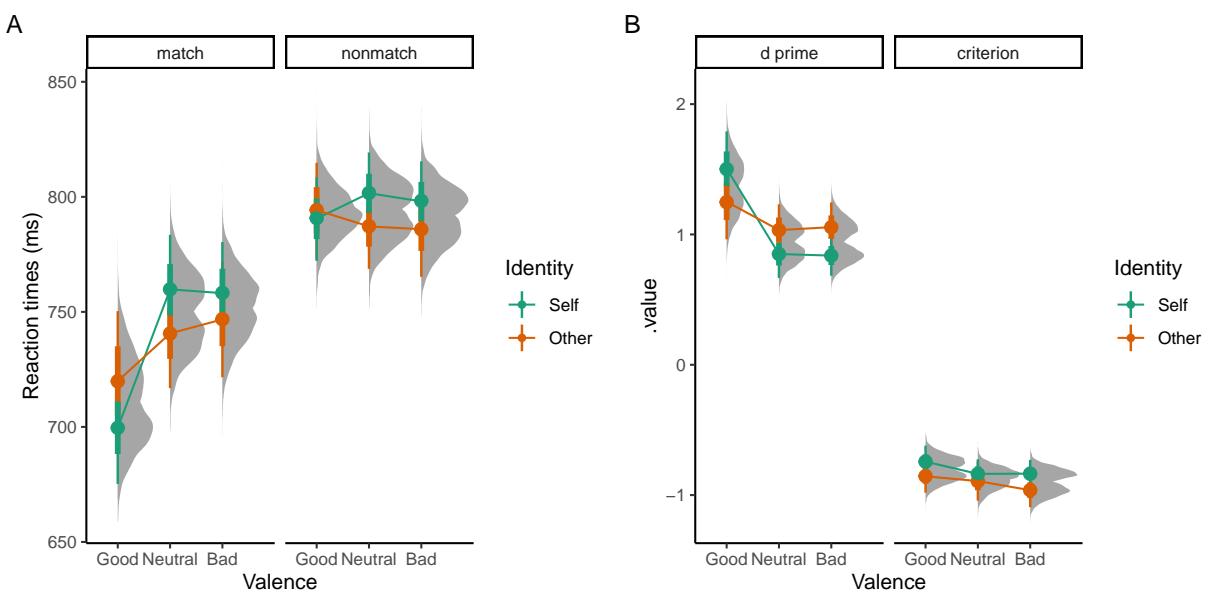


Figure 17. Exp3a: Results of Bayesian GLM analysis.

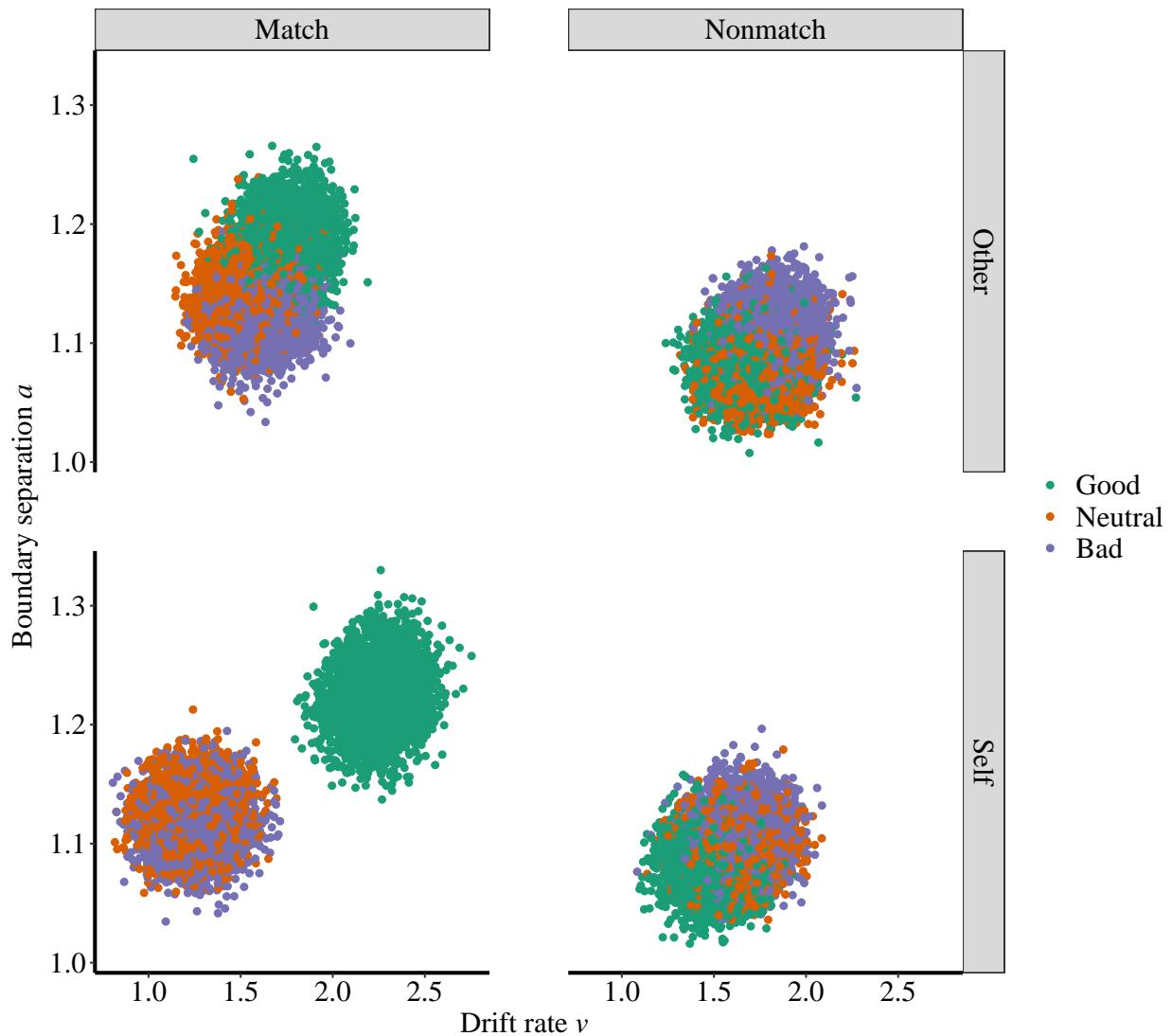


Figure 18. Exp3a: Results of HDDM.

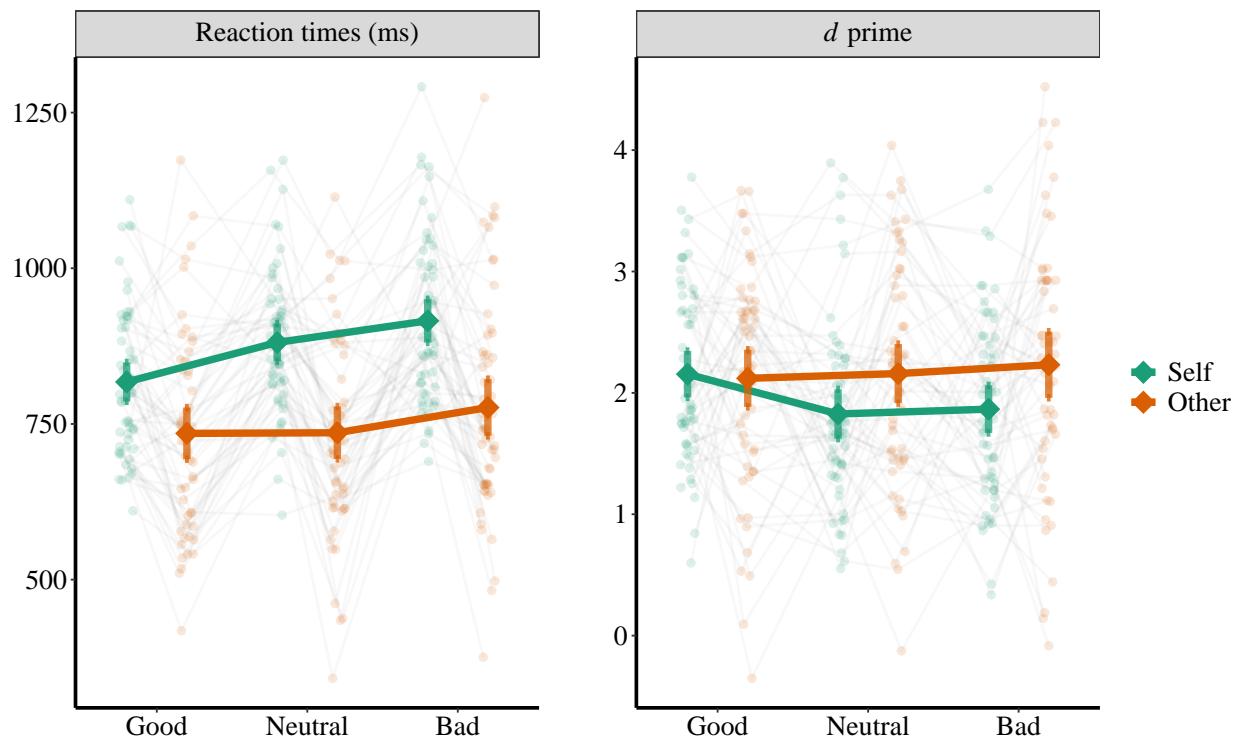


Figure 19. RT and d' of Experiment 3b.

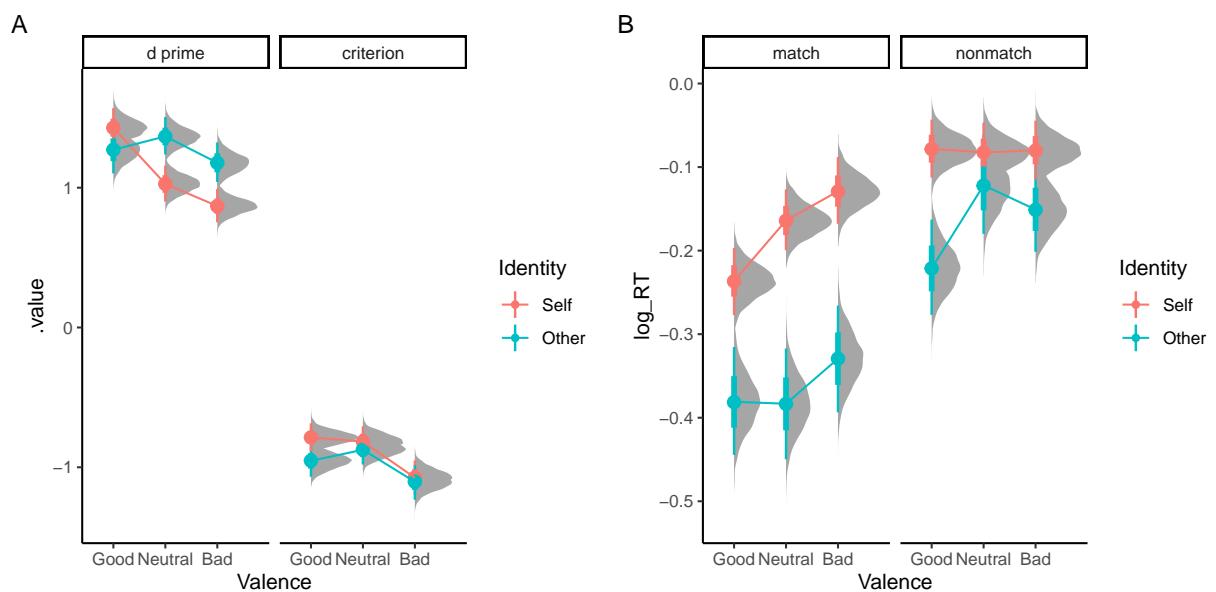


Figure 20. exp3b: Results of Bayesian GLM analysis.

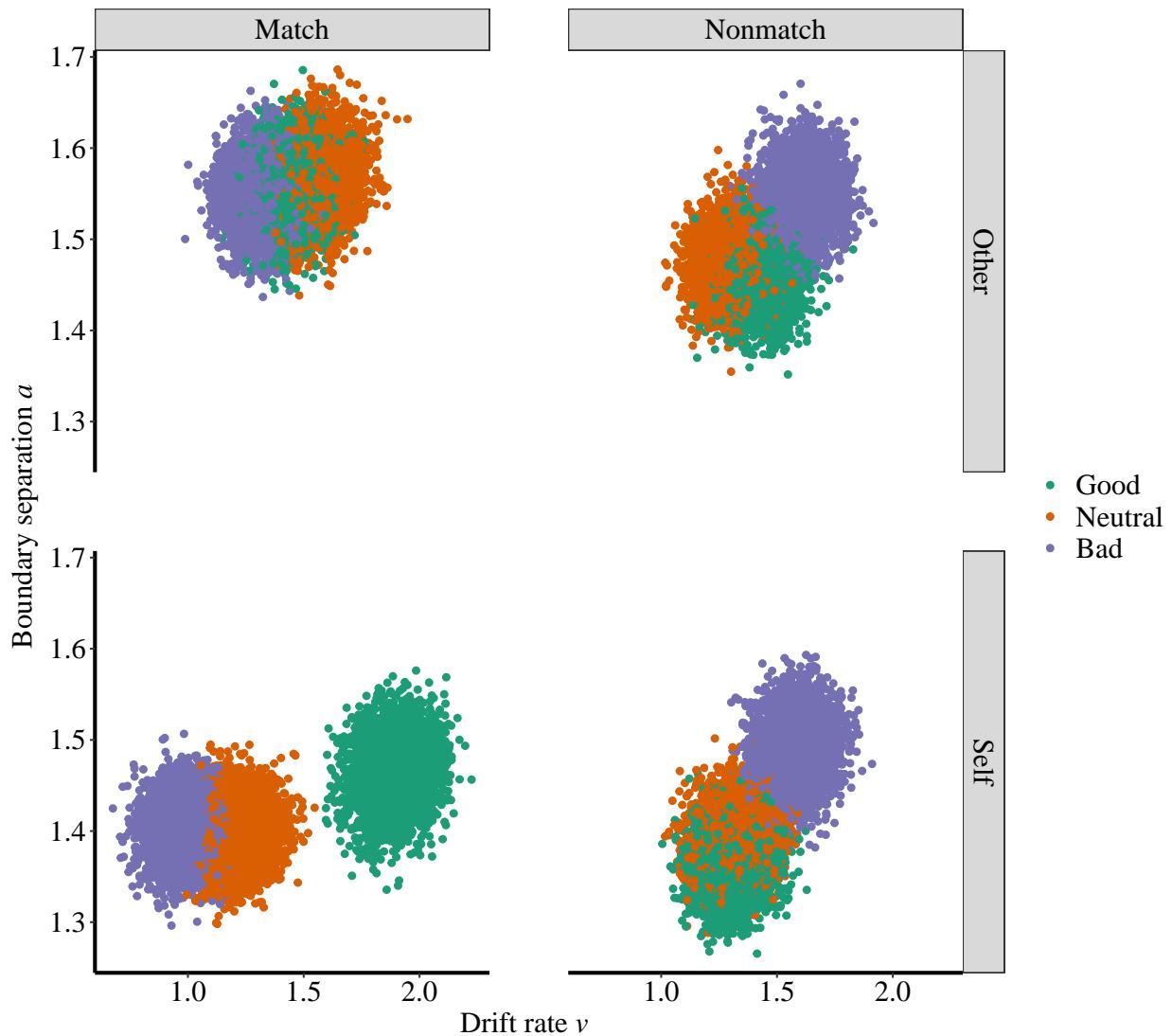


Figure 21. exp3b: Results of HDDM.

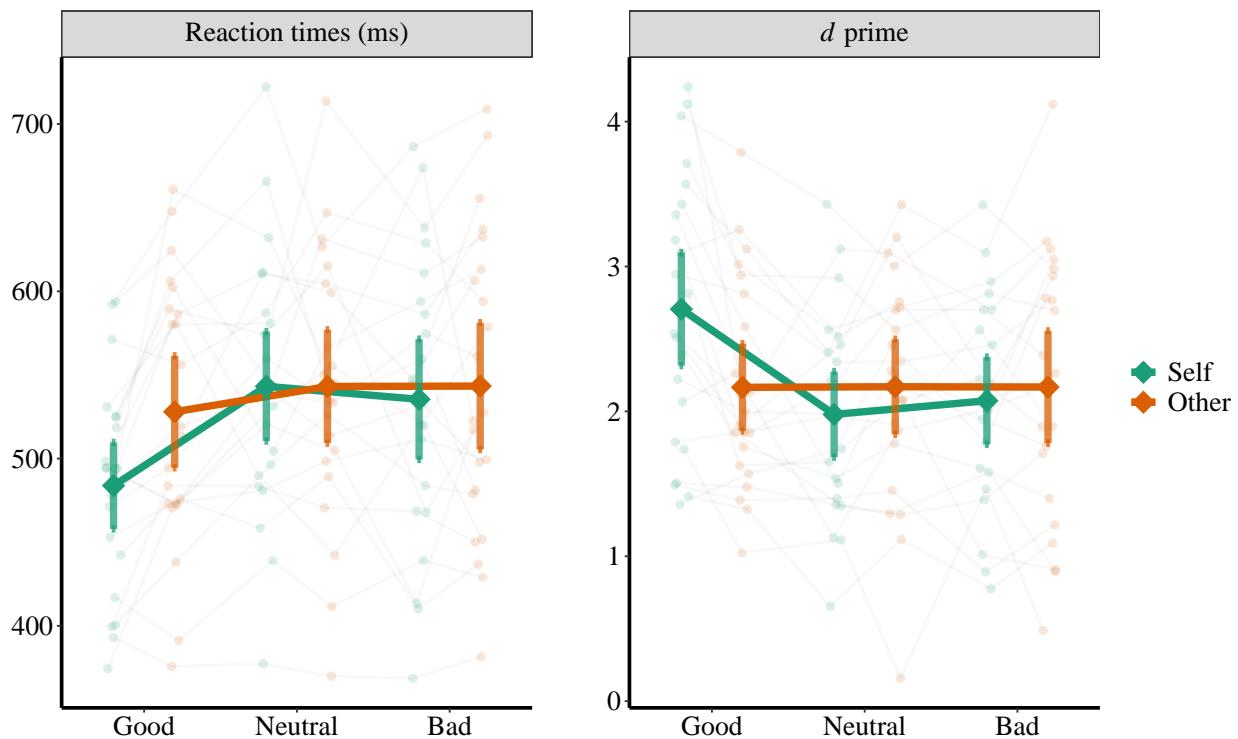


Figure 22. RT and d' of Experiment 6b.

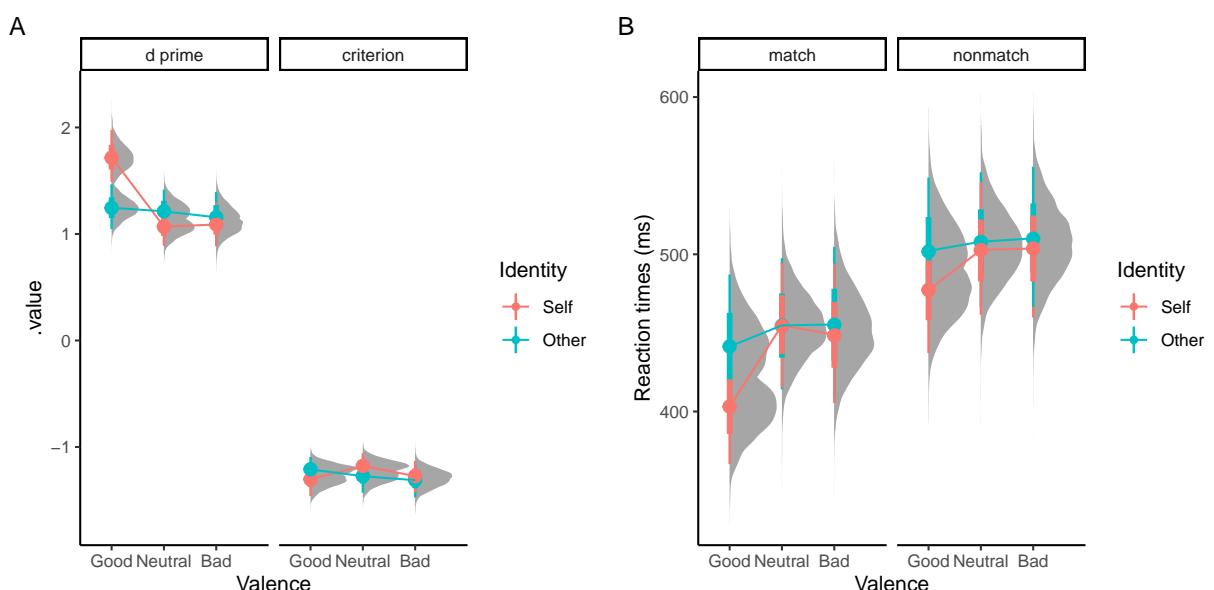


Figure 23. exp6b_d1: Results of Bayesian GLM analysis.

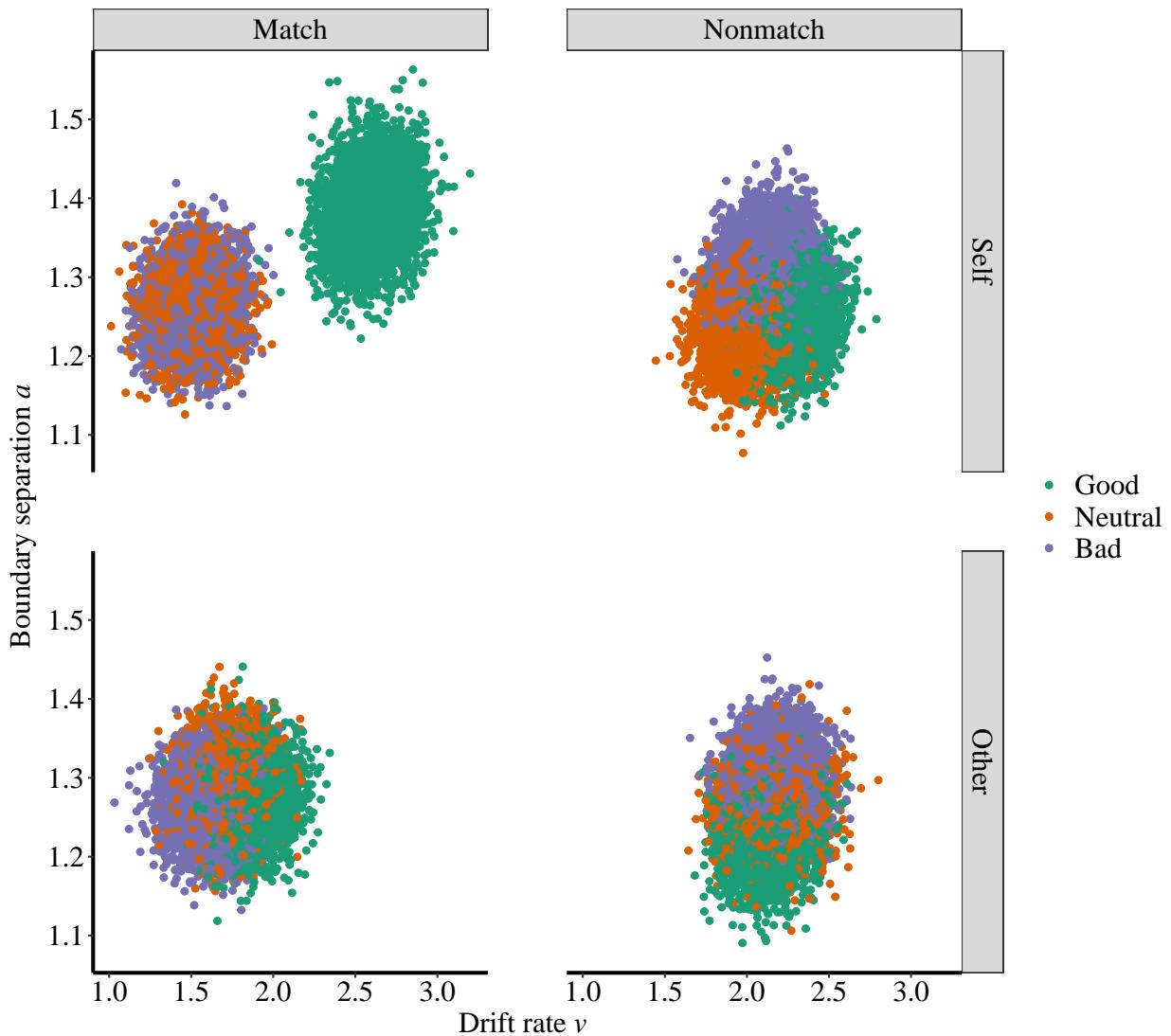


Figure 24. exp6b: Results of HDDM (Day 1).

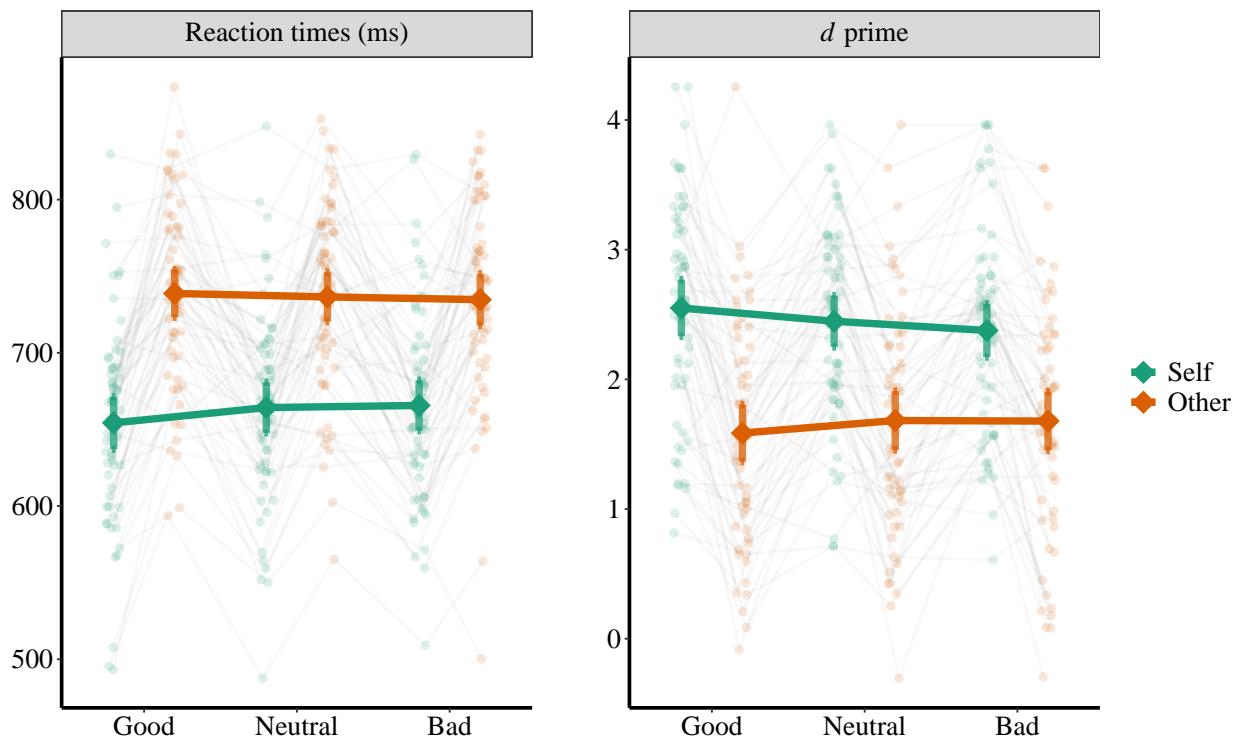


Figure 25. RT and d' of Experiment 4a.

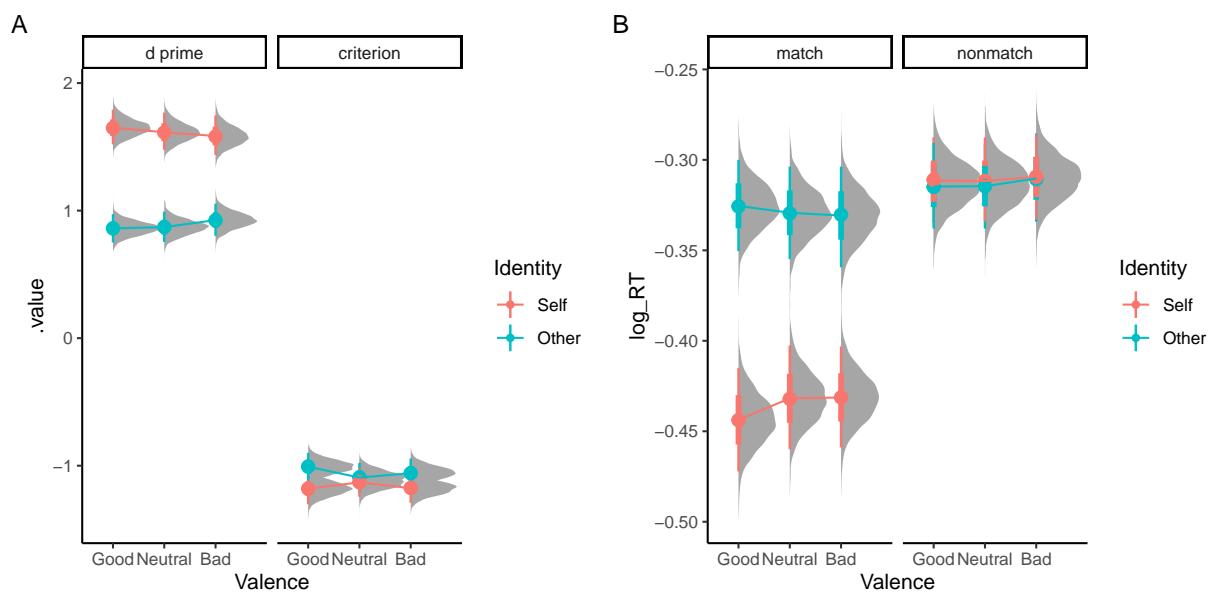


Figure 26. exp4a: Results of Bayesian GLM analysis.

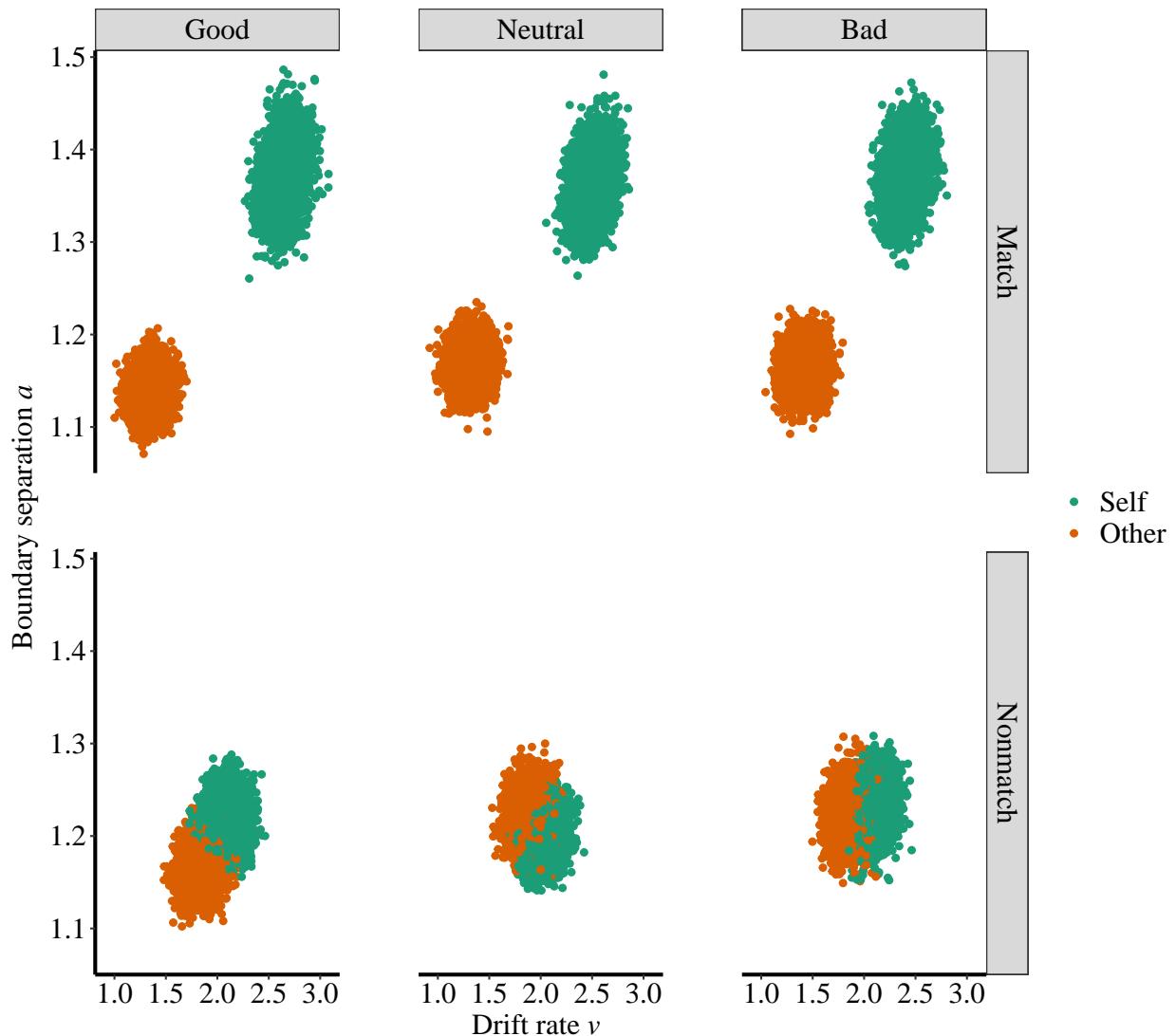


Figure 27. exp4a: Results of HDDM.

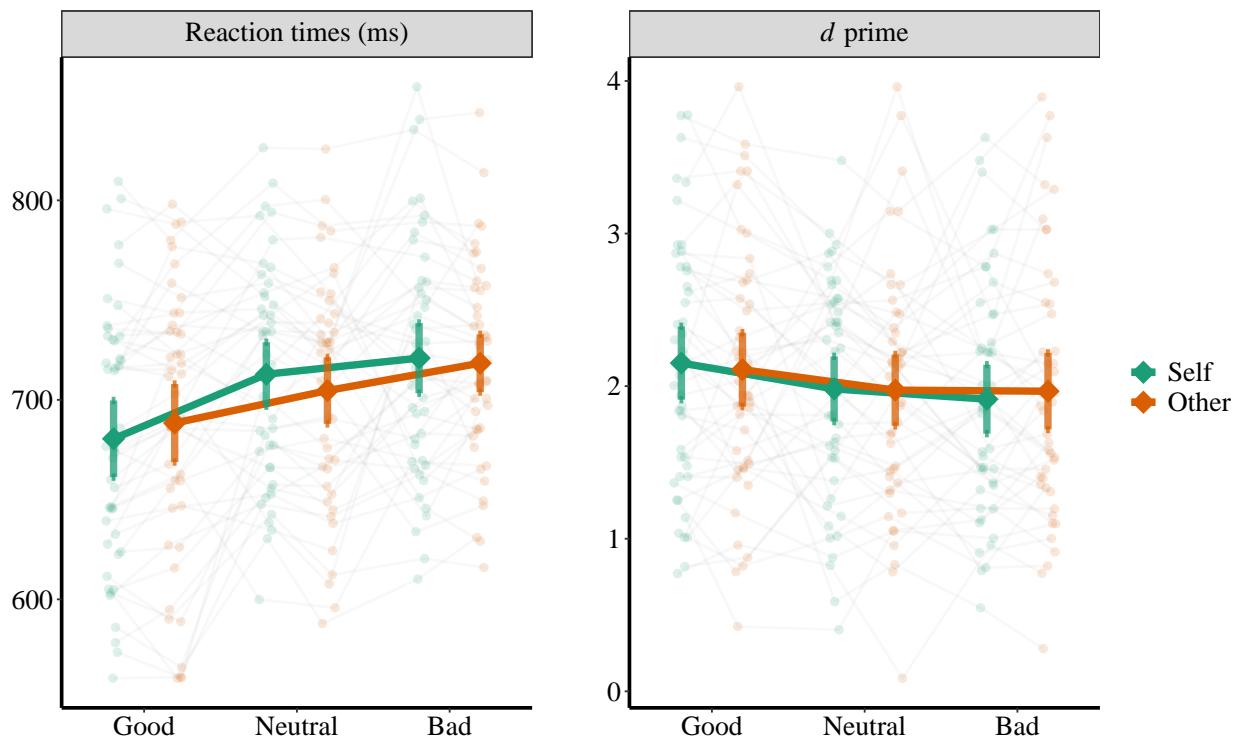


Figure 28. RT and d' prime of Experiment 4b.

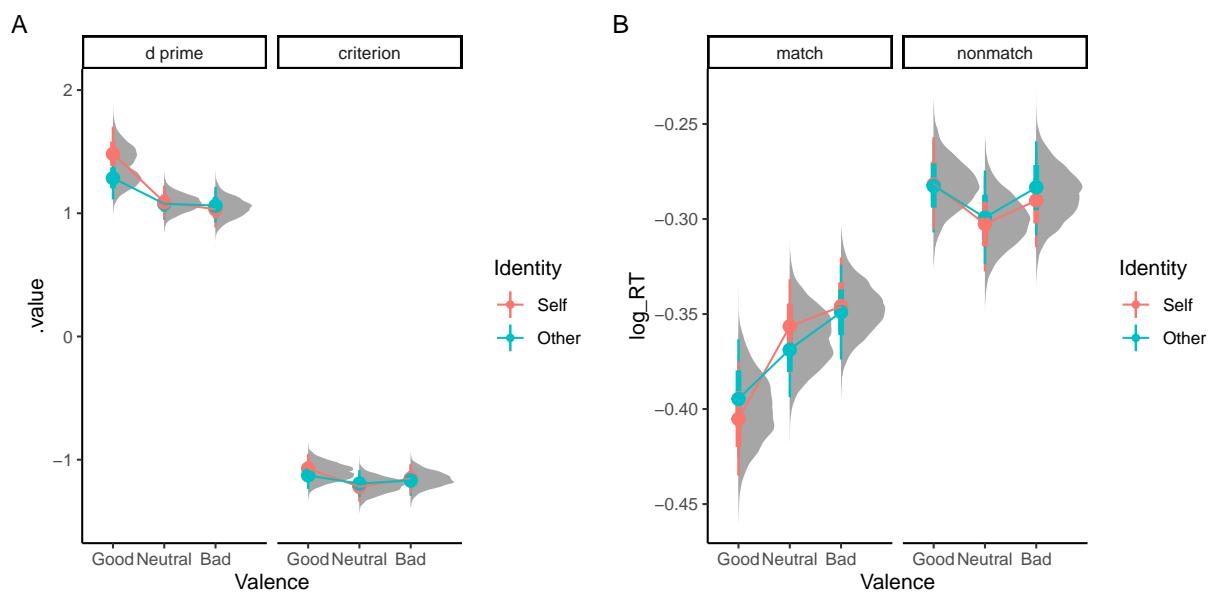


Figure 29. exp4b: Results of Bayesian GLM analysis.

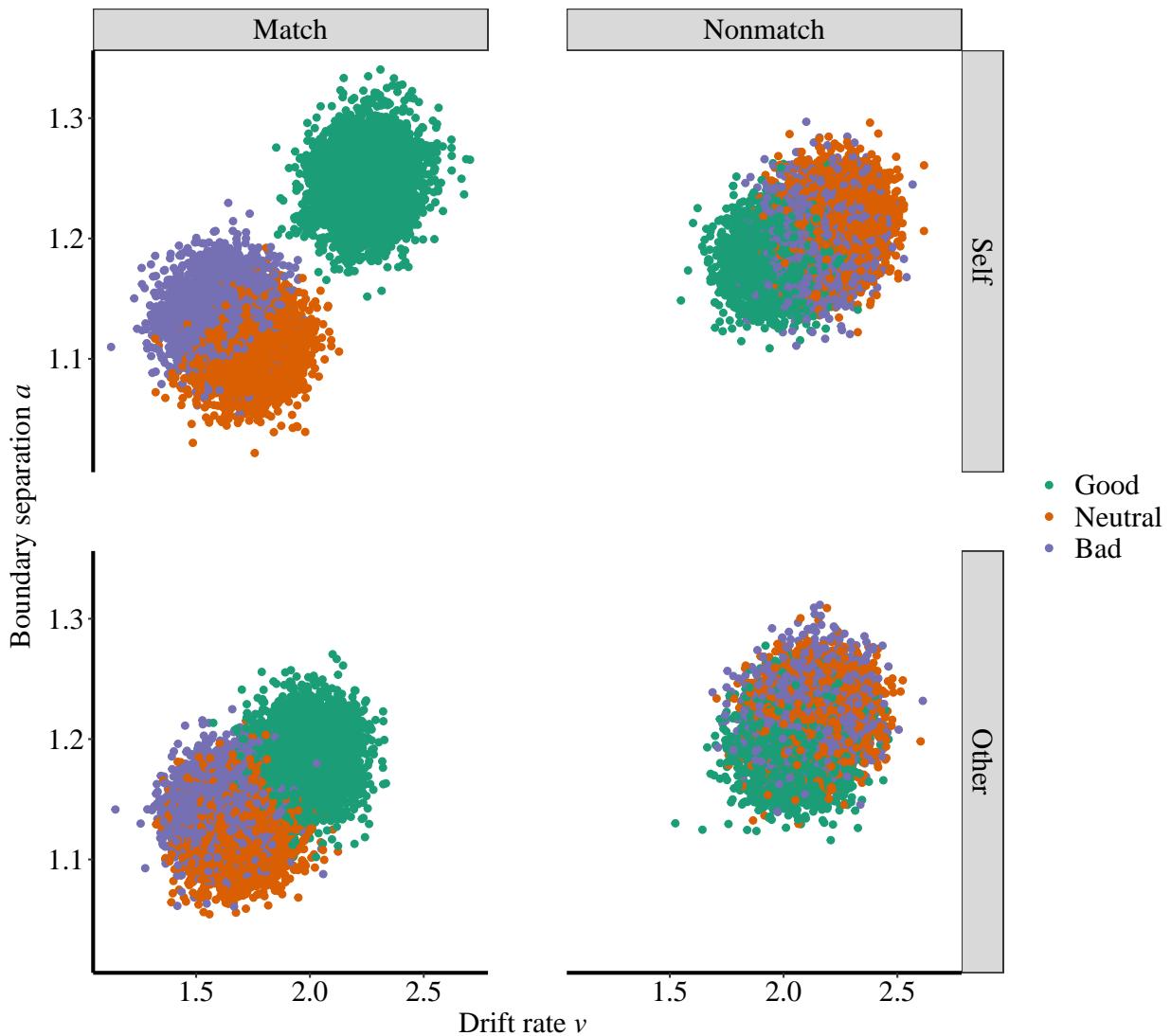


Figure 30. exp4b: Results of HDDM.

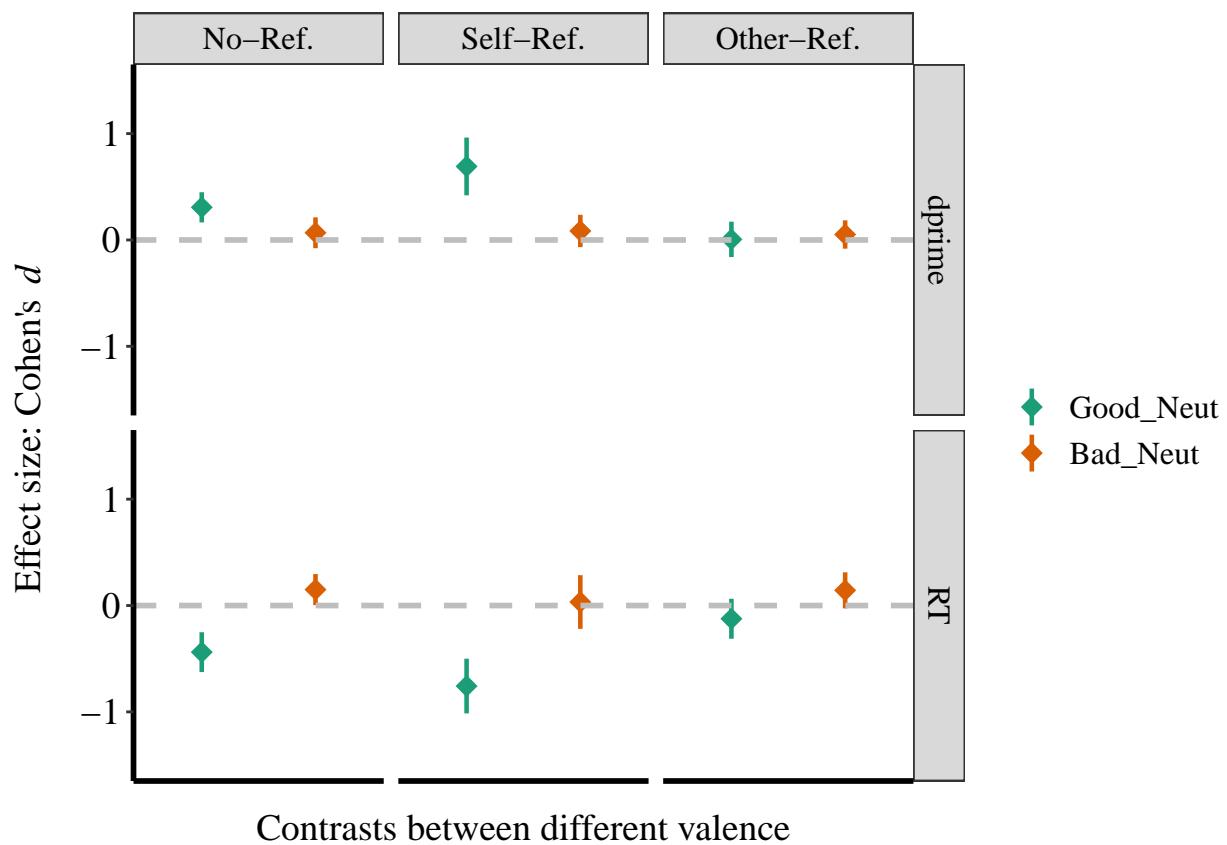


Figure 31. Effect size (Cohen's d) of Valence.

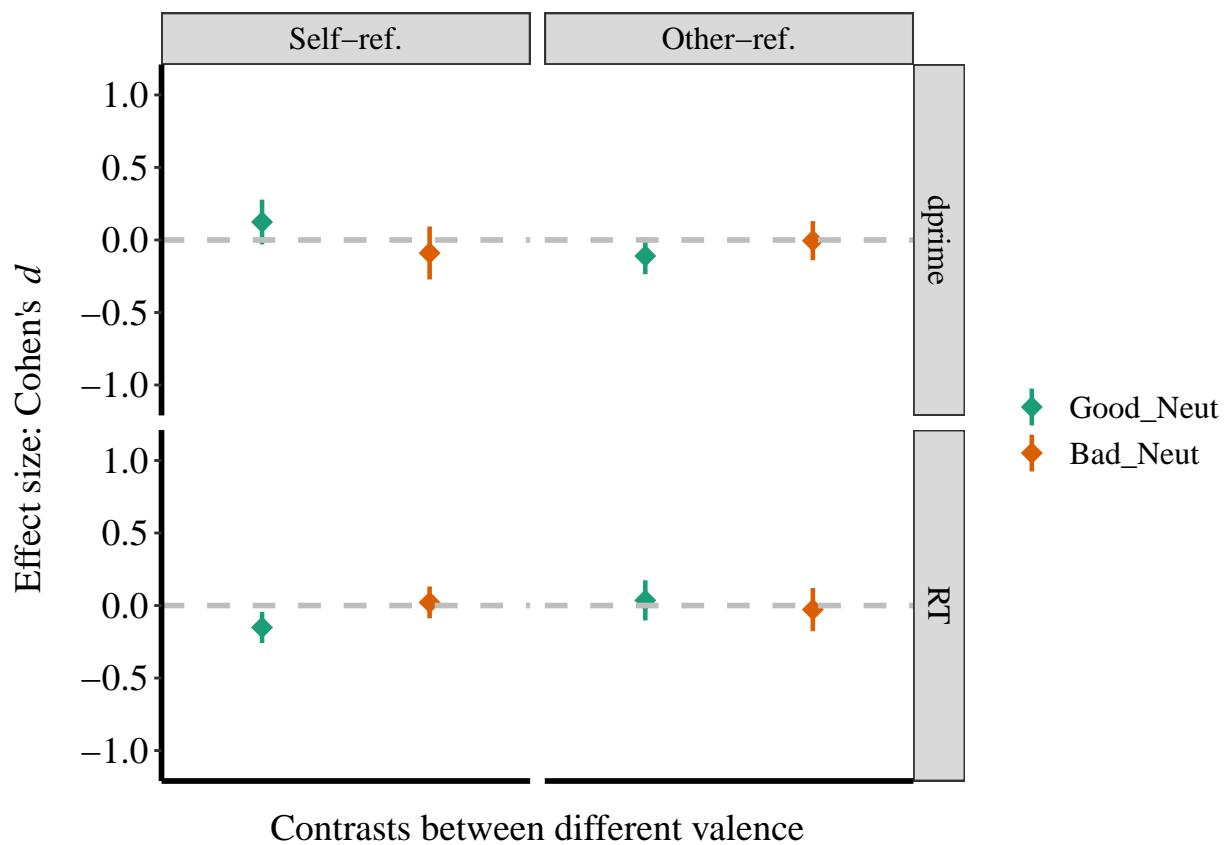


Figure 32. Effect size (Cohen's d) of Valence in Exp4a.

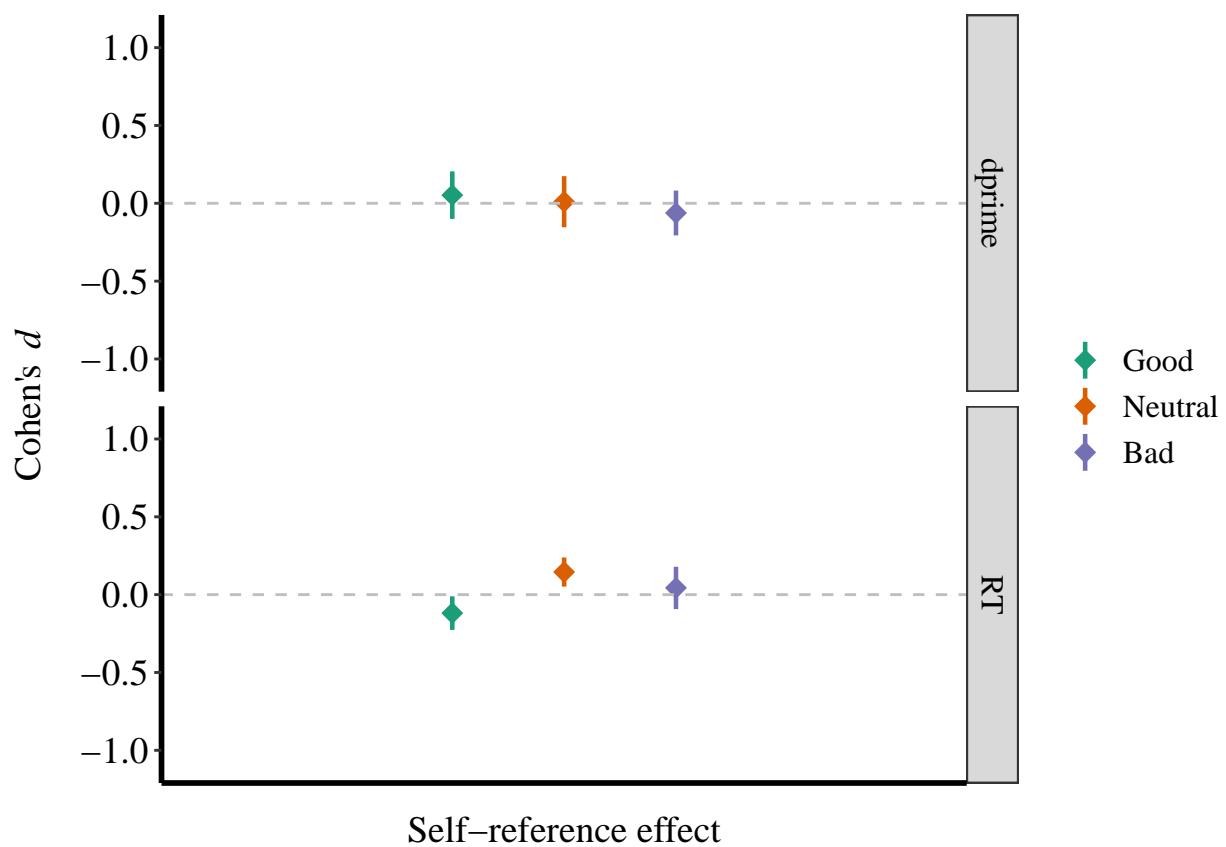


Figure 33. Effect size (Cohen's d) of Valence in Exp4b.

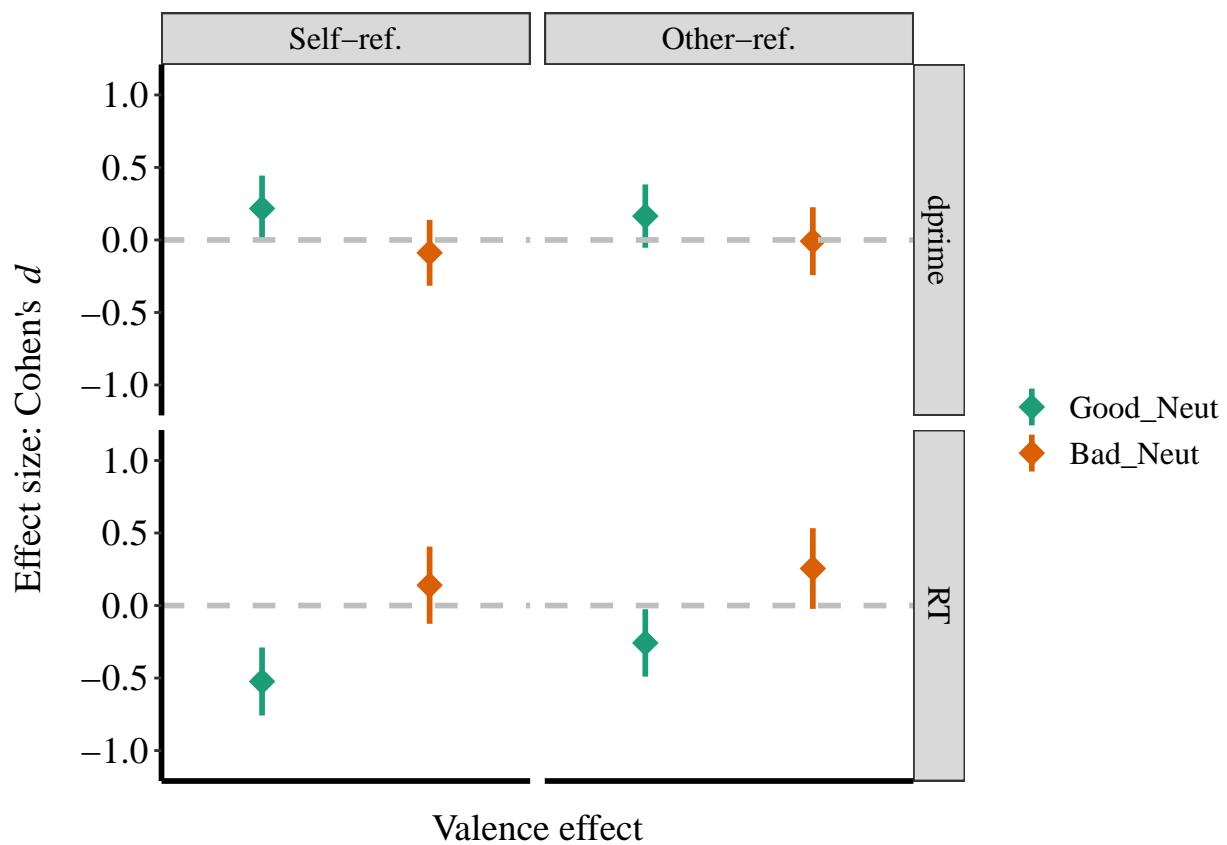


Figure 34. Effect size (Cohen's d) of Valence in Exp4b.

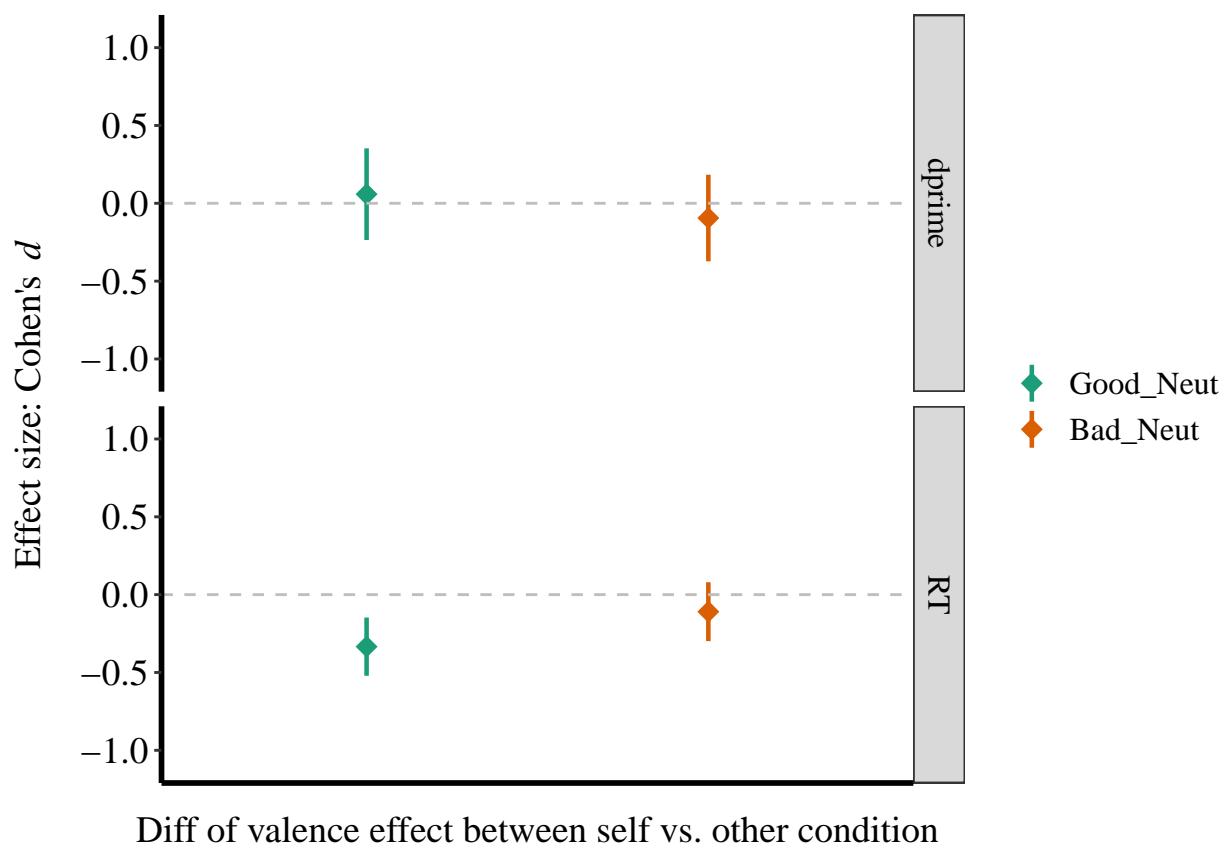


Figure 35. Effect size (Cohen's d) of Valence in Exp4b.

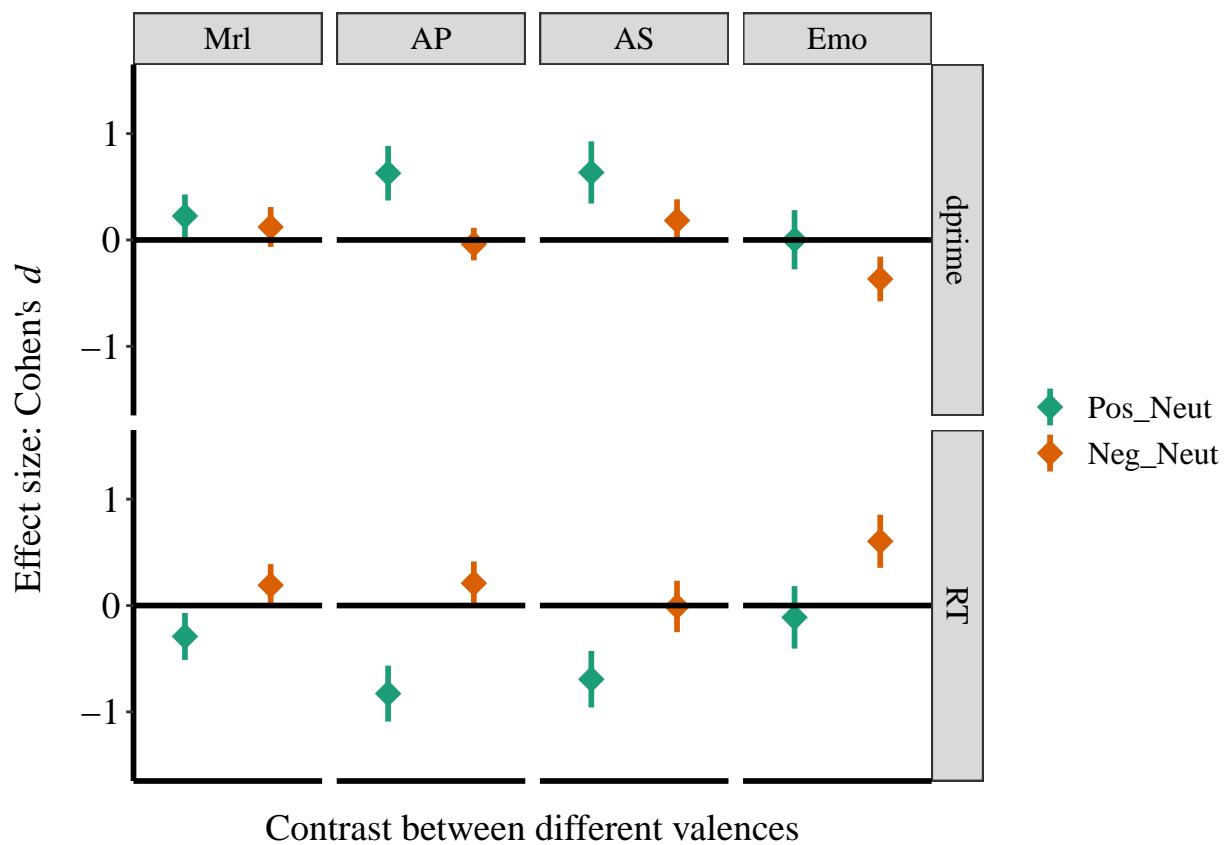


Figure 36. Effect size (Cohen's d) of Valence in Exp5.

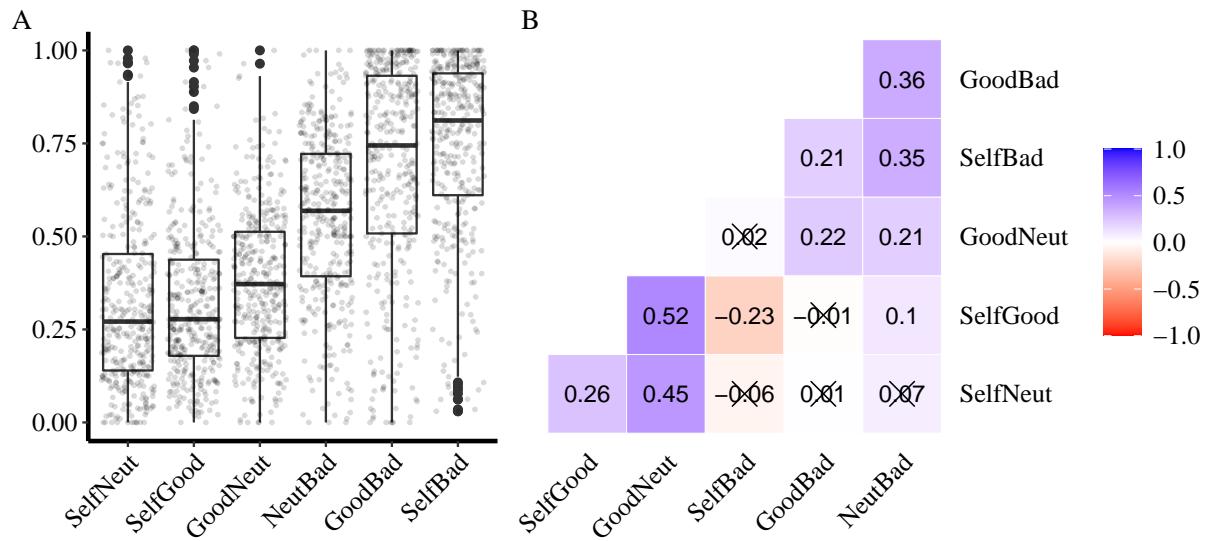


Figure 37. Self-rated personal distance

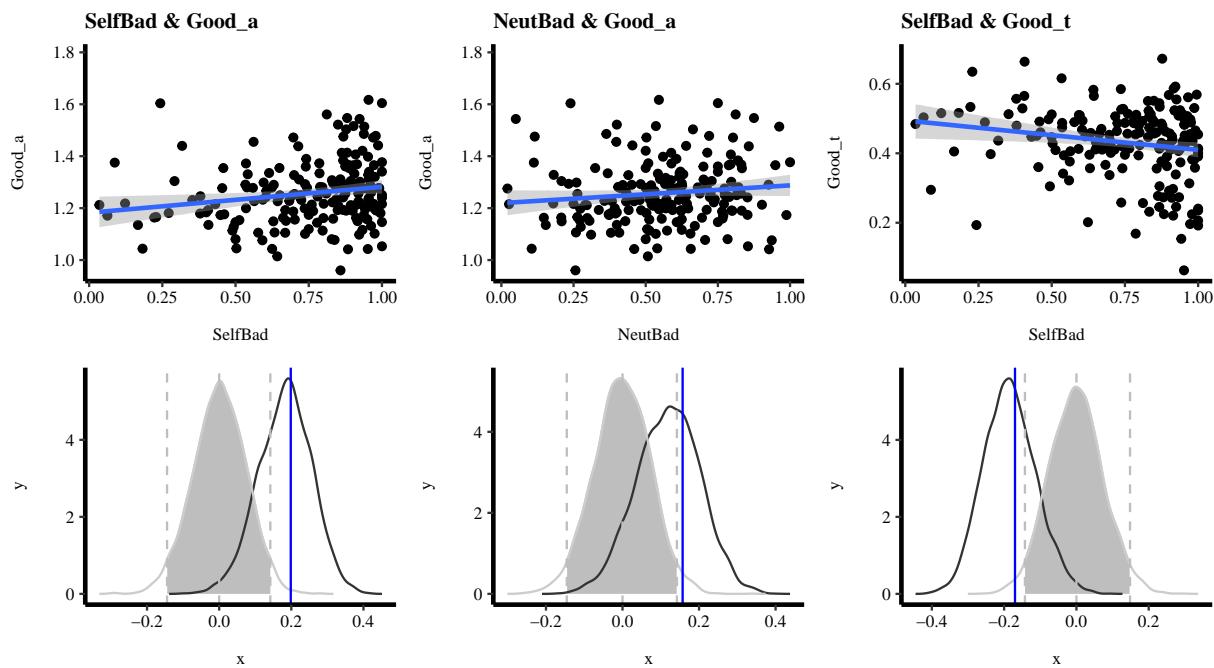


Figure 38. Correlation between moral identity and boundary separation of good condition; moral self-image and drift rate of good condition

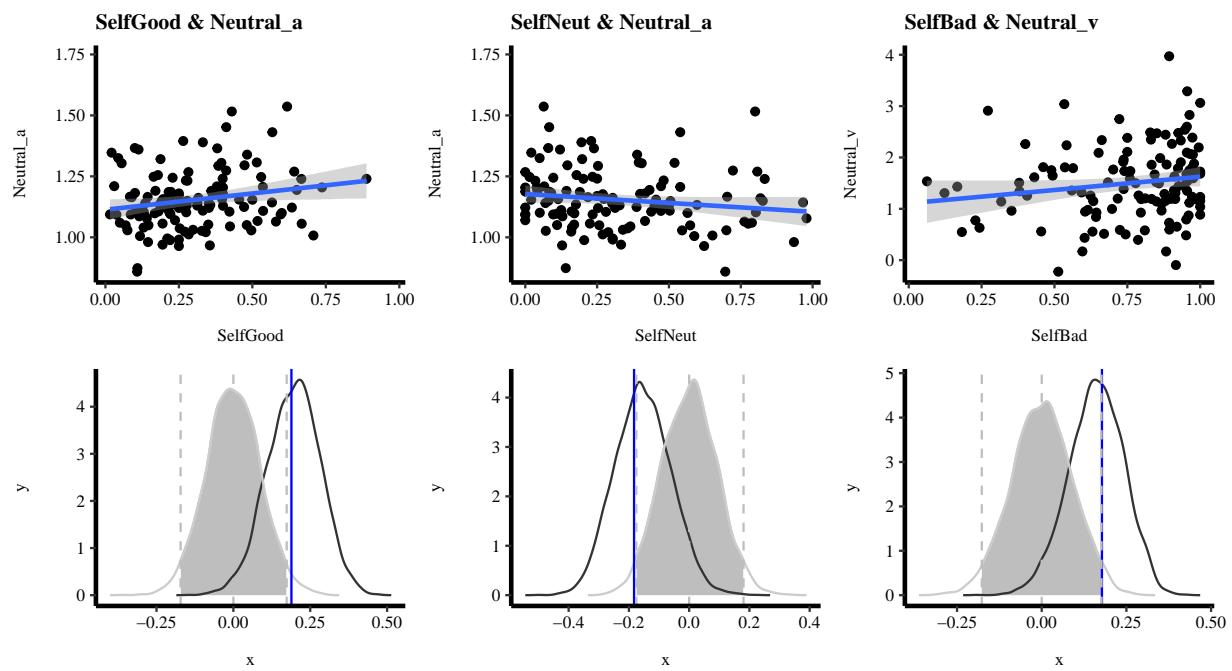


Figure 39. Correlation between personal distance and boundary separation of neutral condition