

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]. Its importance is manifested in many
43 ways in human cognition. For example, morality is a basic dimension of person perception
44 (Abele, Ellemers, Fiske, Koch, & Yzerbyt, 2020; Goodwin, 2015; Goodwin, Piazza, &
45 Rozin, 2014; Willis & Todorov, 2006), moral judgment is common in daily life [].

46 [Given the importance of morality in social life, moral character, i.e., the, is crucial
47 for individual.] More specially, a person needs to both accurately evaluate whether the
48 moral character of others and behave in a way that she is perceived as a moral person, or at
49 least not a morally bad person. The former was usually investigated as person perception
50 in social psychology while the latter was studied separately as moral self-concept, moral
51 self-image, or moral self-enhancement. There are abundant of evidence that people weigh
52 morality heavily in evaluating others (Goodwin, 2015) and evaluating the change of
53 identity of others []. These findings suggest that, given the importance of morality in social
54 life, morality has been internalized in socialized individuals, and this internalized moral
55 concept influence how they perceiving, remembering, and making decisions.

56 When it comes to self perception, there is accumulating evidence that people actively
57 maintain a good moral-self image. For example, recent research found that
58 self-enhancement effect is stronger than that in other domains, such as competence or
59 social competence (Tappin & McKay, 2017). Also, participants maintain their moral
60 self-image even after their own unethical behaviors (e.g., cheating) []. Similarly, when asked
61 how likely they will act ethically or unethically, most participants showed the tendency of
62 less likely to do unethical things []. In other words, existing evidence supported the notion
63 that morality is important in person perception and self-concept, people are motivated to

64 maintain a good moral-self image.

65 [whether moral character information influence perception?, link to exp1a, b, c, and
66 exp2] Yet, as Freeman and Ambady (2011) put it, the focus of the previous studies is not
67 to explain the perceptual process, rather, they are explaining the higher-order social
68 cognitive processes that come after, e.g., impression, evaluation. That is, current studies on
69 person perception is studies without perceptual process. In other words, the perceptual
70 decision-making process of moral character related information is unknown. Without
71 knowledge of these processes work, we can not have a full picture of how moral information
72 is processed in our cognition. As increasing attention paid to perceptual process of social
73 perception, it's clear that perceptual decision-making is strongly influenced by social
74 factors, such as group-categorization, stereotype (Xiao, Coppin, & Bavel, 2016). Given the
75 importance of moral character and that moral character related information has strong
76 influence on learning and memory (Carlson, Maréchal, Oud, Fehr, & Crockett, 2020;
77 Stanley & De Brigard, 2019), one might expect that moral character related information
78 could also play a role in perceptual decision-making.

79 [using associative learning task to study the moral character's influence on perception]
80 Though theoretically possible, no empirical studies had directly addressed this issue. There
81 were only a few studies about the temporal dynamics of judging the trustworthiness of face
82 (e.g., Dzhelyova, Perrett, & Jentzsch, 2012), but trustworthy is not equal to morality. One
83 difficulty of studying moral character's influence on perceptual decision-making is that
84 moral character is a high-level and hidden state instead of observable feature, one also
85 needs moral information, e.g., behavior history, to infer moral character of a person. For
86 example, Anderson et al. (2011) asked participant to first study the behavioral description
87 of faces and then asked participant to perform a perceptual detection task. They found
88 that faces linked with “bad” social behavior are detected faster (but see).

89 An alternative is to use abstract semantic concepts to study how these concepts

90 influence perception. After all, abstract concepts of moral character is part of our daily life
91 and it can be used to identify useful constructs in social life. For example, Also,
92 different levels of information may form a dynamic network in human brain and visual cues
93 can activate more abstract, non-observable personal traits (e.g., aggressiveness) (Amodio,
94 2019; Freeman & Ambady, 2011). If a concept of moral character (e.g., good person) is
95 activated, it should also be able to have influence on perceptual process of the visual cues
96 through the dynamic network, especially when the perceptual decision-making is about the
97 concept-cue association. In this case, abstract concepts of moral character may serve as
98 signal of moral reputation (for others) or moral self-concept. Indeed, previous studies used
99 the moral words and found that moral related information can be perceived faster
100 (Gantman & Van Bavel, 2014). If moral character is an important social category, then, as
101 those other social categories (races, education, etc, see Xiao et al. (2016)), moral character
102 related information might change the perceptual processes.

103 Here, we used an associative learning paradigm to study how moral character concept
104 change perceptual decision-making. In this paradigm, simple geometric shapes were paired
105 with different words whose dominant use is to describe the moral character of person.
106 Participants first learn the association between shapes and words, i.e. good-person and
107 triangle, building direct association between high-level, hidden moral character and visual
108 cue. After remembered the associations, they perform a matching task to judge whether
109 the shape-word pair presented on the screen match the association they learned. This
110 paradigm has been used in studying the perceptual process of self-concept, but had also
111 proven useful in studying other concepts like social group (Enock, Sui, Hewstone, &
112 Humphreys, 2018).

113 Our first question is, whether the words used the in the associative paradigm is really
114 related to the moral character? As we reviewed above, previous theories, especially the
115 interactive dynamic theory, would support this assumption. To validate that moral
116 character concepts activated moral character as a social cue, we used four experiments to

explore and validate the paradigm. The first experiment direct adopted associative paradigm and change the words from “self”, “other” to “good-person”, “neutral-person”, and “bad-person”. Then, we change the words to the ones that have more explicit moral meaning (“kind-person”, “neutral-person”, and “evil-person”). Then, as in Anderson et al. (2011), we asked participant to learn the behavioral history of three different names, and then use the names as moral character words. Finally, we also tested that simultaneously present shape-word pair and sequentially present word and shape didn’t change the pattern. All of these four experiments showed a robust effect of moral content.

[possible explanations: person-based self-categorization vs. stimuli-based valence]

Then, we explored the underlying mechanism. Two theoretical frameworks are relevant here. The first one is valence theory, which emphasize the importance of valence in perceptual decision-making. Under this framework, the valence of the stimuli drives the perceptual decision-making. There exists three different sub-theories under this framework: Negativity effect, positivity effect, and valence effect. The first one might be the threaten detection theory, which predicted that threatening stimuli are preferentially processed because of the evolutionary advantage. Though appealing in the first glance, this threat-detection theory itself has been questioned for the evidence on which this theory was initially proposed. That is negative prioritized because of the low-level physical features of the stimuli []. given that the physical stimuli associated with moral character are manipulated in our study, we expect the low level feature will not play a role, and therefore we have a low a priori on the negative first prediction. The second prediction is positivity hypothesis, which predict that positive moral character will be prioritized. There are also evidence consistent with this idea. For example, XXX found that trustworthy faces attracted attention more than untrustworthy faces, probably because trustworthy faces are more likely to be the collaborative partners subsequent tasks, which will bring reward. A third possibility is that both negative and positive is faster than neutral because of the valence. The underlying of the valence assumption is that the stimuli presented in the

¹⁴⁴ associative task (word-shape pair) can elicit approaching-avoiding motivation. This
¹⁴⁵ probability of this assumption is true is low, because the all these stimuli do not have
¹⁴⁶ visual cue that really threatening or rewarding. Therefore, the difference in perceptual
¹⁴⁷ decision-making make reflect the value of the words to participants, which can be the
¹⁴⁸ interaction of the meaning of the words and the participants' idiosyncratic characteristics.

¹⁴⁹ Previous research converged that, if an object to be of value to an individual, then
¹⁵⁰ that object must be judged as relevant to that individual, i.e., self-relevance (Juechems &
¹⁵¹ Summerfield, 2019; Reicher & Hopkins, 2016). There are two possible way an external
¹⁵² stimuli be valuable (relevant) for an individual. First, we might evaluate its
¹⁵³ rewarding-threatening value, as many previous perceptual research had done. In this
¹⁵⁴ explanation, we will view the moral character and the person behind the moral character,
¹⁵⁵ as objects and only judge whether they are rewarding or potentially rewarding to us. A
¹⁵⁶ different view is that we will still perceive those moral character as human and apply social
¹⁵⁷ categorization, i.e., we categorize whether the person behind the moral character is in the
¹⁵⁸ same group as we do. However, the above four experiments can not distinguish between
¹⁵⁹ these two possibilities, because there are evidence for both reward (Sui, He, & Humphreys,
¹⁶⁰ 2012) and in-group (Enock et al., 2018) prioritization.

¹⁶¹ [Distinguish two explanations by make self salient, and found relative adv of self:
¹⁶² exp3a, 3b, 6b] Though these two framework has similar prediction for the studies with
¹⁶³ moral character such as “good-person”, “neutral person”, and “bad person”, they have
¹⁶⁴ different prediction when if the experiment design include both identity and moral valence.
¹⁶⁵ With an orthogonal design, we have good, bad, and neutral conditions for both self and
¹⁶⁶ other. In this case the identity become salient and participants are less likely to
¹⁶⁷ spontaneously identify good-person as self, but the value of good-person still exists. This
¹⁶⁸ means that the social categorization theory predicts participants prioritize good-self but
¹⁶⁹ not good-other, while reward-based attention theory predicts participants are both
¹⁷⁰ prioritized. Also, as in Hu, Lan, Macrae, and Sui (2020), people may also only identify

¹⁷¹ with good-self instead of neutral or bad self. That is, people will show a unique pattern of
¹⁷² self-identification: only good-self is identified as “self” while all the others categories were
¹⁷³ excluded.

¹⁷⁴ In exp 3a, 3b, and 6b, we found that (1) good-self is always faster than neutral-self
¹⁷⁵ and bad-self, but good-other only have weak to null advantage to neutral-other and
¹⁷⁶ bad-other. which mean the social categorization is self-centered. (2) good-self’s advantage
¹⁷⁷ over good other only occur when self- and other- were in the same task. i.e. the relative
¹⁷⁸ advantage is competition based instead of absolute. These three experiments suggest that
¹⁷⁹ people more like to view the moral character stimuli as person and categorize good-self as
¹⁸⁰ an unique category against all others. A meta-analysis showed that there was no effect of
¹⁸¹ valence when the identity is other.

¹⁸² [what we care? valence of the self exp4a or identity of the good exp4b?] Next, we go
¹⁸³ further to disentangle the self-good complex: people care more about whether the self is
¹⁸⁴ good, or whether the good is self, or both? If people care more about whether the self is
¹⁸⁵ good, then, subtle cue of the valence may have an impact on perceptual process of the self.
¹⁸⁶ In contrast, if people care more about whether the good’s identity, i.e., whether the good is
¹⁸⁷ self, then subtle cue of identity (self v. other) may have a impact on percpetual process of
¹⁸⁸ the good. We tested the good-self complex with two more experiments. In exp 4a (id is
¹⁸⁹ task-relevant, valence is task-irrelevant), if people care about the valence of the self, then,
¹⁹⁰ the task-irrelevant information may influence the processing of the self. While in exp 4b
¹⁹¹ (valence is task-relevant, id is task-irrelevant), if people care about the id of the good, then,
¹⁹² the task-irrelevant id information will has a influence on the process of the good.

¹⁹³ [whether categorize self as positive is not limited to morality? no, but limited to
¹⁹⁴ traits, yet not state] Self-categorization is not limited to morality, even morality is central
¹⁹⁵ to social life. we used aesthetic aspect as another instance.

¹⁹⁶ [Self-bad as an index of self-reported self-categorization in moral term]

197 Below are just my thoughts Still need to thinking about the relationship between
198 several key concepts: self-categorization, morality (moral character), perceptual
199 decision-making, interactive dynamic theory.

200 As previous studies in social psychology and cognitive psychology found the
201 perceiving self and perceiving others has huge difference, therefore, a question immediately
202 followed is, how the self- and other-related moral character information are processed in
203 perceptual decision-making. To investigate this phenomenon, we included both self- and
204 other-related moral character information in the self-tagging paradigm. abstract moral
205 concepts are prioritized in perceptual decision-making? There are several possible
206 predictions.

207 *Self-categories* are cognitive groupings of self and some class of stimuli as identical or
208 different from some other class. [Turner et al.] *Personal identity* refers to self-categories
209 that define the individual as a unique person in terms of his or her individual differences
210 from other (in-group) persons. *Social identity* refers to the shared social categorical self
211 (“us” vs. “them”).

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

215 variable self and morality in perception. self is variable, morality is basic, the point
216 is: how moral other is perceive as extension of self, not reverse.

217 categorization based on morality and variable self/identity

218 What is the prediction of the model/theory?

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

221 People are more likely to identify themselves with trustworthy faces (Verosky &

222 Todorov, 2010) (trustworthy faces has longer RTs).

223 In 1950s, Bruner (1957) had proposed the “New Look” approach of perception, which
224 was resurrected by accumulating evidence (Stolier & Freeman, 2016; Xiao et al., 2016).
225 These studies supported the view that there is a bidirectional interplay between perception
226 and higher-level cognition, such as stereotype (Stolier & Freeman, 2016; Xiao et al., 2016),
227 and self-relevance (Sui et al., 2012). Few studies also tested whether moral-laden
228 information was prioritized in the perception (Anderson et al., 2011; Gantman & Van
229 Bavel, 2014). Still, the moral self-image information has rarely been studied (except Hu et
230 al. (2020)).

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

235 In our experimental setting, we have the following limitations: The perceptual
236 decision-making will show certain pattern under certain task demand. In our case, it's the
237 forced, speed, two-option choice task.

238 Additional assumption about the self-moral id: People will automatically enhance
239 their moral concept, therefore prefer the positive moral self-concept.

240

Disclosures

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

245 All the experiments reported were not pre-registered. Most experiments (1a ~ 6b,
246 except experiment 3b) reported in the current study were first finished between 2014 to

247 2016 in Tsinghua University, Beijing, China. Participants in these experiments were
248 recruited in the local community. To increase the sample size of experiments to 50 or more
249 (Simmons, Nelson, & Simonsohn, 2013), we recruited additional participants in Wenzhou
250 University, Wenzhou, China in 2017 for experiment 1a, 1b, 4a, and 4b. Experiment 3b was
251 finished in Wenzhou University in 2017. To have a better estimation of the effect size, we
252 included the data from two experiments (experiment 7a, 7b) that were reported in Hu et
253 al. (2020) (See Table S1 for overview of these experiments).

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

256 **General methods**

257 **Design and Procedure**

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

270 The experiment 1a, 1b, 1c, 2, and 6a shared a 2 (matching: match vs. nonmatch) by
271 3 (moral valence: good vs. neutral vs. bad) within-subject design. Experiment 1a was the

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

284 For experiment 3a, 3b, 4a, 4b, 6b, 7a, and 7b, we included self-reference as another
285 within-subject variable in the experimental design. For example, the experiment 3a directly
286 extend the design of experiment 1a into a 2 (matchness: match vs. nonmatch) by 2
287 (reference: self vs. other) by 3 (moral valence: good vs. neutral vs. bad) within-subject
288 design. Thus in experiment 3a, there were six conditions (good-self, neutral-self, bad-self,
289 good-other, neutral-other, and bad-other) and six shapes (triangle, square, circle, diamond,
290 pentagon, and trapezoids). The experiment 6b was an EEG experiment extended from
291 experiment 3a but presented the label and shape sequentially. Because of the relatively
292 high working memory load (six label-shape pairs), experiment 6b were conducted in two
293 days: the first day participants finished perceptual matching task as a practice, and the
294 second day, they finished the task again while the EEG signals were recorded. Experiment
295 3b was designed to separate the self-referential trials and other-referential trials. That is,
296 participants finished two different blocks: in the self-referential blocks, they only responded
297 to good-self, neutral-self, and bad-self, with half match trials and half non-match trials; for
298 the other-reference blocks, they only responded to good-other, neutral-other, and

299 bad-other. Experiment 7a and 7b were designed to test the cross task robustness of the
300 effect we observed in the aforementioned experiments (see, Hu et al., 2020). The matching
301 task in these two experiments shared the same design with experiment 3a, but only with
302 two moral valence, i.e., good vs. bad. We didn't include the neutral condition in
303 experiment 7a and 7b because we found that the neutral and bad conditions constantly
304 showed non-significant results in experiment 1 ~ 6.

305 Experiment 4a and 4b were design to test the automaticity of the binding between
306 self/other and moral valence. In 4a, we used only two labels (self vs. other) and two shapes
307 (circle, square). To manipulate the moral valence, we added the moral-related words within
308 the shape and instructed participants to ignore the words in the shape during the task. In
309 4b, we reversed the role of self-reference and valence in the task: participant learnt three
310 labels (good-person, neutral-person, and bad-person) and three shapes (circle, square, and
311 triangle), and the words related to identity, “self” or “other”, were presented in the shapes.
312 As in 4a, participants were told to ignore the words inside the shape during the task.

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

319 E-prime 2.0 was used for presenting stimuli and collecting behavioral responses,
320 except that experiment 7a and 7b used Matlab Psychtoolbox (Brainard, 1997; Pelli, 1997).
321 For participants recruited in Tsinghua University, they finished the experiment individually
322 in a dim-lighted chamber, stimuli were presented on 22-inch CRT monitors and their head
323 were fixed by a chin-rest brace. The distance between participants' eyes and the screen was
324 about 60 cm. The visual angle of geometric shapes was about $3.7^\circ \times 3.7^\circ$, the fixation cross

325 is of ($0.8^\circ \times 0.8^\circ$ of visual angle) at the center of the screen. The words were of $3.6^\circ \times 1.6^\circ$
326 visual angle. The distance between the center of the shape or the word and the fixation
327 cross was 3.5° of visual angle. For participants recruited in Wenzhou University, they
328 finished the experiment in a group consisted of 3 ~ 12 participants in a dim-lighted testing
329 room. Participants were required to finished the whole experiment independently. Also,
330 they were instructed to start the experiment at the same time, so that the distraction
331 between participants were minimized. The stimuli were presented on 19-inch CRT monitor.
332 The visual angles are could not be exactly controlled because participants's chin were not
333 fixed.

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

338 Data analysis

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

343 *Classic NHST.*

344 First, as in Sui et al. (2012), we analyzed the accuracy and reaction times using
345 classic repeated measures ANOVA in the Null Hypothesis Significance Test (NHST)
346 framework. Repeated measures ANOVAs is essentially a two-step mixed model. In the first
347 step, we estimate the parameter on individual level, and in the second step, we used
348 repeated ANOVA to test the Null hypothesis. More specifically, for the accuracy, we used a
349 signal detection approach, in which individual' sensitivity d' was estimated first. To

estimate the sensitivity, we treated the match condition as the signal while the nonmatch conditions as noise. Trials without response were coded either as “miss” (match trials) or “false alarm” (nonmatch trials). Given that the match and nonmatch trials are presented in the same way and had same number of trials across all studies, we assume that participants’ inner distribution of these two types of trials had equal variance but may had different means. That is, we used the equal variance Gaussian SDT model (EVSDT) here (Rouder & Lu, 2005). The d' was then estimated as the difference of the standardized hit and false alarm rats (Stanislaw & Todorov, 1999):

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

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

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

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

Bayesian hierarchical generalized linear model (GLM).

The classic NHST approach may ignore the uncertainty in estimate of the parameters for SDT (Rouder & Lu, 2005), and using mean RT assumes normal distribution of RT data, which is always not true because RTs distribution is skewed (Rousselet & Wilcox, 2019). To better estimate the uncertainty and use a more appropriate model, we also tried

³⁷⁵ Bayesian hierarchical generalized linear model to analyze each experiment's accuracy and
³⁷⁶ RTs data. We used BRMs (Bürkner, 2017) to build the model, which used Stan (Carpenter
³⁷⁷ et al., 2017) to estimate the posterior.

³⁷⁸ In the GLM model, we assume that the accuracy of each trial is Bernoulli distributed
³⁷⁹ (binomial with 1 trial), with probability p_i that $y_i = 1$.

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

³⁸⁰ In the perceptual matching task, the probability p_i can then be modeled as a function of
³⁸¹ the trial type:

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

³⁸² The outcomes y_i are 0 if the participant responded "nonmatch" on trial i , 1 if they
³⁸³ responded "match". The probability of the "match" response for trial i for a participant is
³⁸⁴ p_i . We then write the generalized linear model on the probits (z-scores; Φ , "Phi") of ps . Φ
³⁸⁵ is the cumulative normal density function and maps z scores to probabilities. Given this
³⁸⁶ parameterization, the intercept of the model (β_0) is the standardized false alarm rate
³⁸⁷ (probability of saying 1 when predictor is 0), which we take as our criterion c . The slope of
³⁸⁸ the model (β_1) is the increase of saying 1 when predictor is 1, in z -scores, which is another
³⁸⁹ expression of d' . Therefore, $c = -zHR = -\beta_0$, and $d' = \beta_1$.

³⁹⁰ In each experiment, we had multiple participants, then we need also consider the
³⁹¹ variations between subjects, i.e., a hierarchical mode in which individual's parameter and
³⁹² the population level parameter are estimated simultaneously. We assume that the
³⁹³ outcome of each trial is Bernoulli distributed (binomial with 1 trial), with probability p_{ij}
³⁹⁴ that $y_{ij} = 1$.

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

³⁹⁵ Similarly, the generalized linear model was extended to two levels:

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

³⁹⁶ The outcomes y_{ij} are 0 if participant j responded “nonmatch” on trial i , 1 if they
³⁹⁷ responded “match”. The probability of the “match” response for trial i for subject j is p_{ij} .
³⁹⁸ We again can write the generalized linear model on the probits (z-scores; Φ , “Phi”) of ps .

³⁹⁹ The subjective-specific intercepts ($\beta_0 = -zFAR$) and slopes ($\beta_1 = d'$) are described
⁴⁰⁰ by multivariate normal with means and a covariance matrix for the parameters.

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

⁴⁰¹ For the reaction time, we used the log normal distribution
⁴⁰² ([https://lindeloev.github.io/shiny-rt/#34_\(shifted\)_log-normal](https://lindeloev.github.io/shiny-rt/#34_(shifted)_log-normal)). This distribution has
⁴⁰³ two parameters: μ , σ . μ is the mean of the logNormal distribution, and σ is the disperse of
⁴⁰⁴ the distribution. The log normal distribution can be extended to shifted log normal
⁴⁰⁵ distribution, with one more parameter: shift, which is the earliest possible response.

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

⁴⁰⁶ Shifted log-normal distribution:

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

⁴⁰⁷ y_{ij} is the RT of the i th trial of the j th participants.

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

$$\sigma_j \sim Cauchy()$$

408 *Hierarchical drift diffusion model (HDDM).*

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

419 **Synthesized results.** We also reported the synthesized results from the

420 experiments, because many of them shared the similar experimental design. We reported
 421 the results in five parts: valence effect, explicit interaction between valence and
 422 self-relevance, implicit interaction between valence and self-relevance, specificity of valence
 423 effect, and behavior-questionnaire correlation.

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

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

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

⁴³² M_1 is the mean of the first condition, sd_1 is the standard deviation of the first
⁴³³ condition, while M_2 is the mean of the second condition, sd_2 is the standard deviation of
⁴³⁴ the second condition. r is the correlation coefficient between data from first and second
⁴³⁵ condition. n is the number of data point (in our case the number of participants included
⁴³⁶ in our research).

⁴³⁷ The effect size from each experiment were then synthesized by random effect model
⁴³⁸ using `metafor` (Viechtbauer, 2010). Note that to avoid the cases that some participants
⁴³⁹ participated more than one experiments, we inspected the all available information of
⁴⁴⁰ participants and only included participants' results from their first participation. As
⁴⁴¹ mentioned above, 24 participants were intentionally recruited to participate both exp 1a
⁴⁴² and exp 2, we only included their results from experiment 1a in the meta-analysis.

⁴⁴³ We also estimated the synthesized effect size using Bayesian hierarchical model,
⁴⁴⁴ which extended the two-level hierarchical model in each experiment into three-level model,
⁴⁴⁵ which experiment as an additional level. For SDT, we can use a nested hierarchical model
⁴⁴⁶ to model all the experiment with similar design:

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

⁴⁴⁷ where

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

⁴⁴⁸ The outcomes y_{ijk} are 0 if participant j in experiment k responded “nonmatch” on trial i ,
⁴⁴⁹ 1 if they responded “match”.

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

⁴⁵⁰ and the experiment level parameter mu_{0k} and mu_{1k} is from a higher order
⁴⁵¹ distribution:

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

⁴⁵² in which μ_0 and μ_1 means the population level parameter.

⁴⁵³ This model can be easily expand to three-level model in which participants and
⁴⁵⁴ experiments are two group level variable and participants were nested in the experiments.

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

⁴⁵⁵ y_{ijk} is the RT of the i th trial of the j th participants in the k th experiment.

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

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

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

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

⁴⁵⁶ Using the Bayesian hierarchical model, we can directly estimate the over-all effect of
⁴⁵⁷ valence on d' across all experiments with similar experimental design, instead of using a
⁴⁵⁸ two-step approach where we first estimate the d' for each participant and then use a
⁴⁵⁹ random effect model meta-analysis (Goh et al., 2016).

460 ***Valence effect.***

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

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

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

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

468 In the third part, we focused on experiment 4a and 4b, which were designed to
469 examine the implicit effect of the interaction between moral valence and self-referential
470 processing. We are interested in one particular question: will self-referential and morally
471 positive valence had a mutual facilitation effect. That is, when moral valence (experiment
472 4a) or self-referential (experiment 4a) was presented as task-irrelevant stimuli, whether
473 they would facilitate self-referential or valence effect on perceptual decision-making. For
474 experiment 4a, we reported the comparisons between different valence conditions under the
475 self-referential task and other-referential task. For experiment 4b, we first calculated the
476 effect of valence for both self- and other-referential conditions and then compared the effect
477 size of these three contrast from self-referential condition and from other-referential
478 condition. Note that the results were also analyzed in a standard repeated measure
479 ANOVA (see supplementary materials).

480 ***Specificity of the valence effect.***

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

485 *Behavior-Questionnaire correlation.*

486 Finally, we explored correlation between results from behavioral results and
487 self-reported measures.

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

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

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

504 Then we confirmed the significant results from the data without self- and
505 other-referential (exp1a, 1b, 1c, 2, 5, 6a). In this confirmatory analysis, we used $\alpha =$
506 0.05 and used bootstrap by BootES package (Kirby & Gerlanc, 2013) to estimate the
507 correlation. To avoid false positive, we further determined the threshold for significant by
508 permutation. More specifically, for each pairs that initially with $p < .05$, we randomly
509 shuffle the participants data of each score and calculated the correlation between the
510 shuffled vectors. After repeating this procedure for 5000 times, we choose arrange these

511 5000 correlation coefficients and use the 95% percentile number as our threshold.

512 **Part 1: Moral valence effect**

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

516 **Experiment 1a**

517 **Methods.**

518 ***Participants.***

519 57 college students (38 female, age = 20.75 ± 2.54 years) participated. 39 of them
520 were recruited from Tsinghua University community in 2014; 18 were recruited from
521 Wenzhou University in 2017. All participants were right-handed except one, and all had
522 normal or corrected-to-normal vision. Informed consent was obtained from all participants
523 prior to the experiment according to procedures approved by the local ethics committees. 6
524 participant's data were excluded from analysis because nearly random level of accuracy,
525 leaving 51 participants (34 female, age = 20.72 ± 2.44 years).

526 ***Stimuli and Tasks.***

527 Three geometric shapes were used in this experiment: triangle, square, and circle.

528 These shapes were paired with three labels (bad person, good person or neutral person).
529 The pairs were counterbalanced across participants.

530 ***Procedure.***

531 This experiment had two phases. First, there was a brief learning stage. Participants
532 were asked to learn the relationship between geometric shapes (triangle, square, and circle)
533 and different person (bad person, a good person, or a neutral person). For example, a

534 participant was told, “bad person is a circle; good person is a triangle; and a neutral person
535 is represented by a square.” After participant remember the associations (usually in a few
536 minutes), participants started a practicing phase of matching task which has the exact task
537 as in the experimental task. In the experimental task, participants judged whether
538 shape–label pairs, which were subsequently presented, were correct. Each trial started with
539 the presentation of a central fixation cross for 500 ms. Subsequently, a pairing of a shape
540 and label (good person, bad person, and neutral person) was presented for 100 ms. The
541 pair presented could confirm to the verbal instruction for each pairing given in the training
542 stage, or it could be a recombination of a shape with a different label, with the shape–label
543 pairings being generated at random. The next frame showed a blank for 1100ms.
544 Participants were expected to judge whether the shape was correctly assigned to the person
545 by pressing one of the two response buttons as quickly and accurately as possible within
546 this timeframe (to encourage immediate responding). Feedback (correct or incorrect) was
547 given on the screen for 500 ms at the end of each trial, if no response detected, “too slow”
548 was presented to remind participants to accelerate. Participants were informed of their
549 overall accuracy at the end of each block. The practice phase finished and the experimental
550 task began after the overall performance of accuracy during practice phase achieved 60%.
551 For participants from the Tsinghua community, they completed 6 experimental blocks of 60
552 trials. Thus, there were 60 trials in each condition (bad-person match, bad-person
553 nonmatch, good-person match, good-person nonmatch, neutral-person match, and
554 neutral-person nonmatch). For the participants from Wenzhou University, they finished 6
555 blocks of 120 trials, therefore, 120 trials for each condition.

556 ***Data analysis.***

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

560 **Results.**

561 ***Classic NHST.***

562 *d prime.*

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

565 We found the effect of Valence ($F(1.96, 97.84) = 6.19, MSE = 0.27, p = .003,$
 566 $\hat{\eta}_G^2 = .020$). The post-hoc comparison with multiple comparison correction revealed that
 567 the shapes associated with Good-person (2.11, SE = 0.14) has greater *d prime* than shapes
 568 associated with Bad-person (1.75, SE = 0.14), $t(50) = 3.304, p = 0.0049$. The Good-person
 569 condition was also greater than the Neutral-person condition (1.95, SE = 0.16), but didn't
 570 reach statistical significant, $t(50) = 1.54, p = 0.28$. Neither the Neutral-person condition is
 571 significantly greater than the Bad-person condition, $t(50) = 2.109, p = .098$.

572 *Reaction times.*

573 We conducted 2 (Matchness: match v. nonmatch) by 3 (Valence: good, neutral, bad)
 574 repeated measure ANOVA. We found the main effect of Matchness ($F(1, 50) = 232.39,$
 575 $MSE = 948.92, p < .001, \hat{\eta}_G^2 = .104$), main effect of valence ($F(1.87, 93.31) = 9.62,$
 576 $MSE = 1,673.86, p < .001, \hat{\eta}_G^2 = .016$), and interaction between Matchness and Valence
 577 ($F(1.73, 86.65) = 8.52, MSE = 1,441.75, p = .001, \hat{\eta}_G^2 = .011$).

578 We then carried out two separate ANOVA for Match and Mismatched trials. For
 579 matched trials, we found the effect of valence . We further examined the effect of valence
 580 for both self and other for matched trials. We found that shapes associated with Good
 581 Person (684 ms, SE = 11.5) responded faster than Neutral (709 ms, SE = 11.5), $t(50) =$
 582 -2.265, $p = 0.0702$ and Bad Person (728 ms, SE = 11.7), $t(50) = -4.41, p = 0.0002$), and
 583 the Neutral condition was faster than the Bad condition, $t(50) = -2.495, p = 0.0415$). For
 584 non-matched trials, there was no significant effect of Valence ()�.

585 ***Bayesian hierarchical GLM.***

586 *d prime.*

587 We fitted a Bayesian hierarchical GLM for signal detection theory approach. The
 588 results showed that when the shapes were tagged with labels with different moral valence,
 589 the sensitivity (d') and criteria (c) were both influence. For the d' , we found that the
 590 shapes tagged with morally good person (2.46, 95% CI[2.21 2.72]) is greater than shapes
 591 tagged with moral bad (2.07, 95% CI[1.83 2.32]), $P_{PosteriorComparison} = 1$. Shape tagged
 592 with morally good person is also greater than shapes tagged with neutral person (2.23,
 593 95% CI[1.95 2.49]), $P_{PosteriorComparison} = 0.97$. Also, the shapes tagged with neutral
 594 person is greater than shapes tagged with morally bad person, $P_{PosteriorComparison} = 0.92$.

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

600 *Reaction times.*

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

609 **HDDM.**

610 We fitted our data with HDDM, using the response-coding (See also, Hu et al., 2020).
 611 We estimated separate drift rate (v), non-decision time (T_0), and boundary separation (a)

for each condition. We found that the shapes tagged with good person has higher drift rate and higher boundary separation than shapes tagged with both neutral and bad person. Also, the shapes tagged with neutral person has a higher drift rate than shapes tagged with bad person, but not for the boundary separation. Finally, we found that shapes tagged with bad person had longer non-decision time (see Figure 3).

Experiment 1b

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

Method.

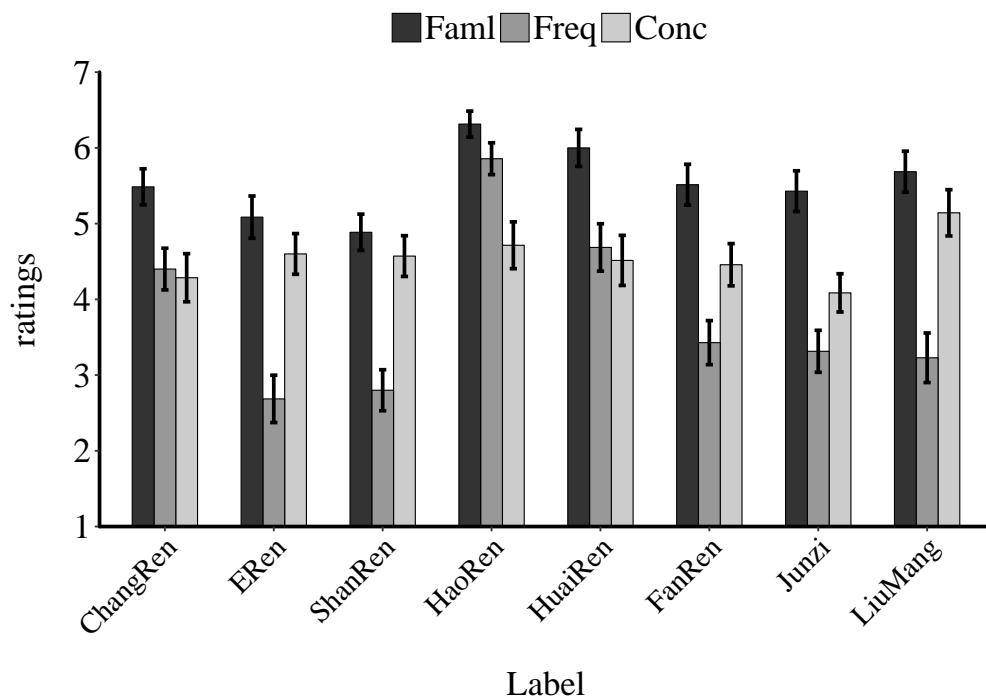
Participants.

72 college students (49 female, age = 20.17 ± 2.08 years) participated. 39 of them were recruited from Tsinghua University community in 2014; 33 were recruited from Wenzhou University in 2017. All participants were right-handed except one, 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 the local ethics committees. 20 participant's data were excluded from analysis because nearly random level of accuracy, leaving 52 participants (36 female, age = 20.25 ± 2.31 years).

Stimuli and Tasks. Three geometric shapes (triangle, square, and circle, with $3.7^\circ \times 3.7^\circ$ of visual angle) were presented above a white fixation cross subtending $0.8^\circ \times 0.8^\circ$ of visual angle at the center of the screen. The three shapes were randomly assigned to three labels with different moral valence: a morally bad person (" ", ERen), a morally good person (" ", ShanRen) or a morally neutral person (" ", ChangRen). The order of the associations between shapes and labels was counterbalanced across participants. Three labels used in this experiment is selected based on the rating results from an independent survey, in which participants rated the familiarity, frequency, and concreteness of eight

637 different words online. Of the eight words, three of them are morally positive (HaoRen,
 638 ShanRen, Junzi), two of them are morally neutral (ChangRen, FanRen), and three of them
 639 are morally negative (HuaiRen, ERen, LiuMang). An independent sample consist of 35
 640 participants (22 females, age 20.6 ± 3.11) were recruited to rate these words. Based on the
 641 ratings (see supplementary materials Figure S1), we selected ShanRen, ChangRen, and
 642 ERen to represent morally positive, neutral, and negative person.

Ratings for each label



643

Procedure.

644 For participants from both Tsinghua community and Wenzhou community, the
 645 procedure in the current study was exactly same as in experiment 1a.
 646

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

Results.

NHST.

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

649 d prime.

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

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

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

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

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

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

Reaction times.

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

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

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

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

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

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

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

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

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

BGLM.

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 influence. For the d' , we found that the shapes tagged with morally good

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

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

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

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

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

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

683 *Reaction time.*

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

692 **HDDM.**

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

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

703 **Experiment 1c**

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

705 Instead of using label with moral valence, we used valence-neutral names in China.

706 Participant first learn behaviors of the different person, then, they associate the names and

707 shapes. And then they perform a name-shape matching task.

708 **Method.**

709 ***Participants.***

710 23 college students (15 female, age = 22.61 ± 2.62 years) participated. All of them

711 were recruited from Tsinghua University community in 2014. Informed consent was

712 obtained from all participants prior to the experiment according to procedures approved by

713 the local ethics committees. No participant was excluded because they overall accuracy

714 were above 0.6.

715 ***Stimuli and Tasks.***

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

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

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

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

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

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

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

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

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

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

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

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

⁷²⁸ SD = 1.1), $t(33) = 30.78$, $p < .001$. Neutral condition was also significant higher than bad
⁷²⁹ conditions $t(33) = 13.9$, $p < .001$ (See supplementary materials).

⁷³⁰ **Procedure.**

⁷³¹ After arriving the lab, participants were informed to complete two experimental
⁷³² tasks, first a social memory task to remember three person and their behaviors, after tested
⁷³³ for their memory, they will finish a perceptual matching task. In the social memory task,
⁷³⁴ the descriptions of three person were presented without time limitation. Participant
⁷³⁵ self-paced to memorized the behaviors of each person. After they memorizing, a
⁷³⁶ recognition task was used to test their memory effect. Each participant was required to
⁷³⁷ have over 95% accuracy before preceding to matching task. The perceptual learning task
⁷³⁸ was followed, three names were randomly paired with geometric shapes. Participants were
⁷³⁹ required to learn the association and perform a practicing task before they start the formal
⁷⁴⁰ experimental blocks. They kept practicing until they reached 70% accuracy. Then, they
⁷⁴¹ would start the perceptual matching task as in experiment 1a. They finished 6 blocks of
⁷⁴² perceptual matching trials, each have 120 trials.

⁷⁴³ **Data Analysis.** Data was analyzed as in experiment 1a.

⁷⁴⁴ **Results.** Figure 7 shows d prime and reaction times of experiment 1c. We
⁷⁴⁵ conducted same analysis as in Experiment 1a. Our analysis didn't show effect of valence
⁷⁴⁶ on d prime, $F(1.93, 42.56) = 0.23$, $MSE = 0.41$, $p = .791$, $\hat{\eta}_G^2 = .005$. Neither the effect of
⁷⁴⁷ valence on RT ($F(1.63, 35.81) = 0.22$, $MSE = 2,212.71$, $p = .761$, $\hat{\eta}_G^2 = .001$) or
⁷⁴⁸ interaction between valence and matchness on RT ($F(1.79, 39.43) = 1.20$,
⁷⁴⁹ $MSE = 1,973.91$, $p = .308$, $\hat{\eta}_G^2 = .005$).

⁷⁵⁰ **Signal detection theory analysis of accuracy.**

⁷⁵¹ We fitted a Bayesian hierarchical GLM for SDT. The results showed that when the
⁷⁵² shapes were tagged with labels with different moral valence, the sensitivity (d') and criteria
⁷⁵³ (c) were both influenced. For the d' , we found that the shapes tagged with morally good

754 person (2.30, 95% CI[1.93 2.70]) is greater than shapes tagged with moral bad (2.11, 95%
 755 CI[1.83 2.42]), $P_{PosteriorComparison} = 0.8$. Shape tagged with morally good person is also
 756 greater than shapes tagged with neutral person (2.16, 95% CI[1.88 2.45]),
 757 $P_{PosteriorComparison} = 0.75$.

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

763 ***Reaction time.***

764 We fitted a Bayesian hierarchical GLM for RTs, with a log-normal distribution as the
 765 link function. We used the posterior distribution of the regression coefficient to make
 766 statistical inferences. As in previous studies, the matched conditions are much faster than
 767 the mismatched trials ($P_{PosteriorComparison} = .75$). We focused on matched trials only, and
 768 compared different conditions: Good () is not faster than the neutral (),
 769 $P_{PosteriorComparison} = .5$, it was faster than the Bad condition (),
 770 $P_{PosteriorComparison} = .88$. And the neutral condition is faster than the bad condition,
 771 $P_{PosteriorComparison} = .95$. However, the mismatched trials are largely overlapped.

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

779 Experiment 2: Sequential presenting

780 Experiment 2 was conducted for two purpose: (1) to further confirm the facilitation
781 effect of positive moral associations; (2) to test the effect of expectation of occurrence of
782 each pair. In this experiment, after participant learned the association between labels and
783 shapes, they were presented a label first and then a shape, they then asked to judge
784 whether the shape matched the label or not (see (Sui, Sun, Peng, & Humphreys, 2014)).
785 Previous studies showed that when the labels presented before the shapes, participants
786 formed expectations about the shape, and therefore a top-down process were introduced
787 into the perceptual matching processing. If the facilitation effect of positive moral valence
788 we found in experiment 1 was mainly drive by top-down processes, this sequential
789 presenting paradigm may eliminate or attenuate this effect; if, however, the facilitation
790 effect occurred because of button-up processes, then, similar facilitation effect will appear
791 even with sequential presenting paradigm.

792 Method.**793 Participants.**

794 35 participants (17 female, age = 21.66 ± 3.03) were recruited. 24 of them had
795 participated in Experiment 1a (9 male, mean age = 21.9, s.d. = 2.9), and the time gap
796 between these experiment 1a and experiment 2 is at least six weeks. The results of 1
797 participants were excluded from analysis because of less than 60% overall accuracy,
798 remains 34 participants (17 female, age = 21.74 ± 3.04).

799 Procedure.

800 In Experiment 2, the sequential presenting makes the matching task much easier than
801 experiment 1. To avoid ceiling effect on behavioral data, we did a few pilot experiments to
802 get optimal parameters, i.e., the conditions under which participant have similar accuracy
803 as in Experiment 1 (around 70 ~ 80% accuracy). In the final procedure, the label (good
804 person, bad person, or neutral person) was presented for 50 ms and then masked by a

805 scrambled image for 200 ms. A geometric shape followed the scrambled mask for 50 ms in
806 a noisy background (which was produced by first decomposing a square with $\frac{3}{4}$ gray area
807 and $\frac{1}{4}$ white area to small squares with a size of 2×2 pixels and then re-combine these
808 small pieces randomly), instead of pure gray background in Experiment 1. After that, a
809 blank screen was presented 1100 ms, during which participants should press a button to
810 indicate the label and the shape match the original association or not. Feedback was given,
811 as in study 1. The next trial then started after 700 ~ 1100 ms blank. Other aspects of
812 study 2 were identical to study 1.

813 ***Data analysis.***

814 Data was analyzed as in study 1a.

815 **Results.**

816 **NHST.**

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

819 *d prime.*

820 There was evidence for the main effect of valence, $F(1.83, 60.36) = 14.41$,
821 $MSE = 0.23$, $p < .001$, $\hat{\eta}_G^2 = .066$. Paired t test showed that the Good-Person condition
822 (2.79 ± 0.17) was with greater d prime than Netural condition (2.21 ± 0.16 , $t(33) = 4.723$,
823 $p = 0.001$) and Bad-person condition (2.41 ± 0.14), $t(33) = 4.067$, $p = 0.008$). There was
824 no-significant difference between Neutral-person and Bad-person conidition, $t(33) = -1.802$,
825 $p = 0.185$.

826 *Reaction time.*

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

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

830 $MSE = 605.36, p < .001, \hat{\eta}_G^2 = .017$) and then analyzed the matched trials and
 831 mismatched trials separately, as in experiment 1a. For matched trials, we found the effect
 832 of valence $F(1.99, 65.76) = 10.57, MSE = 1,192.65, p < .001, \hat{\eta}_G^2 = .067$. Post-hoc t -tests
 833 revealed that shapes associated with Good Person (548 ± 9.4) were responded faster than
 834 Neutral-Person (582 ± 10.9), ($t(33) = -3.95, p = 0.0011$) and Bad Person (582 ± 10.2),
 835 $t(33) = -3.9, p = 0.0013$). While there was no significant differences between Neutral and
 836 Bad-Person condition ($t(33) = -0.01, p = 0.999$). For non-matched trials, there was no
 837 significant effect of Valence ($F(1.99, 65.83) = 0.17, MSE = 489.80, p = .843, \hat{\eta}_G^2 = .001$).

838 **BGLMM.**

839 *Signal detection theory analysis of accuracy.*

840 We fitted a Bayesian hierarchical GLM for SDT. The results showed that when the
 841 shapes were tagged with labels with different moral valence, the sensitivity (d') and criteria
 842 (c) were both influence. For the d' , we found that the shapes tagged with morally good
 843 person ($2.46, 95\% \text{ CI}[2.21 2.72]$) is greater than shapes tagged with moral bad ($2.07, 95\%$
 844 $\text{CI}[1.83 2.32]$), $P_{PosteriorComparison} = 1$. Shape tagged with morally good person is also
 845 greater than shapes tagged with neutral person ($2.23, 95\% \text{ CI}[1.95 2.49]$),
 846 $P_{PosteriorComparison} = 0.97$. Also, the shapes tagged with neutral person is greater than
 847 shapes tagged with morally bad person, $P_{PosteriorComparison} = 0.92$.

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

853 *Reaction times.*

854 We fitted a Bayesian hierarchical GLM for RTs, with a log-normal distribution as the
 855 link function. We used the posterior distribution of the regression coefficient to make

856 statistical inferences. As in previous studies, the matched conditions are much faster than
857 the mismatched trials ($P_{PosteriorComparison} = .75$). We focused on matched trials only, and
858 compared different conditions: Good () is not faster than the neutral (),
859 $P_{PosteriorComparison} = .5$, it was faster than the Bad condition (),
860 $P_{PosteriorComparison} = .88$. And the neutral condition is faster than the bad condition,
861 $P_{PosteriorComparison} = .95$. However, the mismatched trials are largely overlapped.

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

870 Discussion

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

877 Experiment 6a: EEG study 1

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

Method.***Participants.***

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

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

Stimuli.

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

Procedure.

The procedure was similar to Experiment 2. Participants finished 9 blocks of trial,

906 each with 120 trials. In total, participants finished 180 trials for each combination of
907 condition.

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

916 All the stimuli were presented on a gray background (RGB: 127, 127, 127). E-primed
917 2.0 was used to present stimuli and collect behavioral results. Data were collected and
918 analyzed when accuracy performance in total reached 60%.

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

920 **Results.**

921 **NHST.**

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

924 *d* prime.

925 We conducted repeated measures ANOVA, with moral valence as independent
926 variable. The results revealed the main effect of valence ($F(1.74, 40.05) = 3.76$,
927 $MSE = 0.10$, $p = .037$, $\hat{\eta}_G^2 = .021$). Post-hoc analysis revealed that shapes link with Good
928 person (mean = 3.13, SE = 0.109) is greater than Neutral condition (mean = 2.88, SE =
929 0.14), $t = 2.916$, $df = 24$, $p = 0.02$, p-value adjusted by Tukey method, but the *d* prime
930 between Good and bad (mean = 3.03, SE = 0.142) ($t = 1.512$, $df = 24$, $p = 0.3034$, p-value

931 adjusted by Tukey method), bad and neutral ($t = 1.599$, $df = 24$, $p = 0.2655$, p-value
932 adjusted by Tukey method) were not significant.

933 *Reaction times.*

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

936 We found interaction between Matchness and Valence ($F(1.97, 45.20) = 20.45$,
937 $MSE = 450.47$, $p < .001$, $\hat{\eta}_G^2 = .021$) and then analyzed the matched trials and
938 mismatched trials separately, as in experiment 2. For matched trials, we found the effect of
939 valence $F(1.97, 45.25) = 32.37$, $MSE = 522.42$, $p < .001$, $\hat{\eta}_G^2 = .078$. For non-matched
940 trials, there was no significant effect of Valence ($F(1.77, 40.67) = 0.35$, $MSE = 242.15$,
941 $p = .679$, $\hat{\eta}_G^2 = .000$). Post-hoc t -tests revealed that shapes associated with Good Person
942 (mean = 550, SE = 13.8) were responded faster than Neutral-Person (501, SE = 14.7),
943 ($t(24) = -5.171$, $p = 0.0001$) and Bad Person (523, SE = 16.3), $t(24) = -8.137$, $p <$
944 0.0001), and Neutral is faster than Bad-Person condition ($t(32) = -3.282$, $p = 0.0085$).

945 **BGLM.**

946 *Signal detection theory analysis of accuracy.*

947 *Reaction time.*

948 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et
949 al., 2020). We estimated separate drift rate (v), non-decision time (T_0), and boundary
950 separation (a) for each condition. We found that, similar to experiment 2, the shapes
951 tagged with good person has higher drift rate and higher boundary separation than shapes
952 tagged with both neutral and bad person, but only for the self-referential condition. Also,
953 the shapes tagged with neutral person has a higher drift rate than shapes tagged with bad
954 person, but not for the boundary separation, and this effect also exist only for the
955 self-referential condition.

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

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

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

964 **Experiment 3a**

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

969 **Method.**

970 ***Participants.***

971 38 college students (15 female, age = 21.92 ± 2.16) participated in experiment 3a.
972 All of them were right-handed, and all had normal or corrected-to-normal vision. Informed
973 consent was obtained from all participants prior to the experiment according to procedures
974 approved by a local ethics committee. One female and one male student did not finish the
975 experiment, and 1 participants' data were excluded from analysis because less than 60%
976 overall accuracy, remains 35 participants (13 female, age = 22.11 ± 2.13).

977 ***Design.***

978 Study 3a combined moral valence with self-relevance, hence the experiment has a $2 \times$
979 3×2 within-subject design. The first variable was self-relevance, include two levels:

980 self-relevance vs. stranger-relevance; the second variable was moral valence, include good,
981 neutral and bad; the third variable was the matching between shape and label: match
982 vs. nonmatch.

983 ***Stimuli.***

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

989 ***Procedure.***

990 After being fully explained and signed the informed consent, participants were
991 instructed to chose a name that can represent a stranger with same gender as the
992 participant themselves, from a common Chinese name pool. Before experiment, the
993 experimenter explained the meaning of each label to participants. For example, the “good
994 self” mean the morally good side of themselves, them could imagine the moment when they
995 do something’s morally applauded, “bad self” means the morally bad side of themselves,
996 they could also imagine the moment when they doing something morally wrong, and
997 “neutral self” means the aspect of self that does not related to morality, they could imagine
998 the moment when they doing something irrelevant to morality. In the same sense, the
999 “good other”, “bad other”, and “neutral other” means the three different aspects of the
1000 stranger, whose name was chosen before the experiment. Then, the experiment proceeded
1001 as study 1a. Each participant finished 6 blocks, each have 120 trials. The sequence of trials
1002 was pseudo-randomized so that there are 10 matched trials for each condition and 10
1003 non-matched trials for each condition (good self, neutral self, bad self, good other, neutral
1004 other, bad other) for each block.

1005 ***Data Analysis.***

1006 Data analysis followed strategies described in the general method section. Reaction
1007 times and d prime data were analyzed as in study 1 and study 2, except that one more
1008 within-subject variable (i.e., self-relevance) was included in the analysis.

1009 **Results.**

1010 **NHST.**

1011 Figure 16 shows d prime and reaction times of experiment 3a. Less than 5% correct
1012 trials with less than 200ms reaction times were excluded.

1013 *d prime.*

1014 There was evidence for the main effect of valence, $F(1.89, 64.37) = 11.09$,
1015 $MSE = 0.23$, $p < .001$, $\hat{\eta}_G^2 = .039$, and main effect of self-relevance, $F(1, 34) = 3.22$,
1016 $MSE = 0.54$, $p = .082$, $\hat{\eta}_G^2 = .015$, as well as the interaction, $F(1.79, 60.79) = 3.39$,
1017 $MSE = 0.43$, $p = .045$, $\hat{\eta}_G^2 = .022$.

1018 We then conducted separated ANOVA for self-referential and other-referential trials.
1019 The valence effect was shown for the self-referential conditions, $F(1.65, 56.25) = 13.98$,
1020 $MSE = 0.31$, $p < .001$, $\hat{\eta}_G^2 = .119$. Post-hoc test revealed that the Good-Self condition
1021 (1.97 ± 0.14) was with greater d prime than Neutral condition (1.41 ± 0.12 , $t(34) = 4.505$,
1022 $p = 0.0002$), and Bad-self condition (1.43 ± 0.102), $t(34) = 3.856$, $p = 0.0014$. There was
1023 difference between neutral and bad condition, $t(34) = -0.238$, $p = 0.9694$. However, no
1024 effect of valence was found for the other-referential condition $F(1.98, 67.36) = 0.38$,
1025 $MSE = 0.35$, $p = .681$, $\hat{\eta}_G^2 = .004$.

1026 *Reaction time.*

1027 We found interaction between Matchness and Valence ($F(1.98, 67.44) = 26.29$,
1028 $MSE = 730.09$, $p < .001$, $\hat{\eta}_G^2 = .025$) and then analyzed the matched trials and nonmatch
1029 trials separately, as in previous experiments.

1030 For the match trials, we found that the interaction between identity and valence,

1031 $F(1.72, 58.61) = 3.89$, $MSE = 2,750.19$, $p = .032$, $\hat{\eta}_G^2 = .019$, as well as the main effect of
 1032 valence $F(1.98, 67.34) = 35.76$, $MSE = 1,127.25$, $p < .001$, $\hat{\eta}_G^2 = .079$, but not the effect of
 1033 identity $F(1, 34) = 0.20$, $MSE = 3,507.14$, $p = .660$, $\hat{\eta}_G^2 = .001$. As for the d prime, we
 1034 separated analyzed the self-referential and other-referential trials. For the Self-referential
 1035 trials, we found the main effect of valence, $F(1.80, 61.09) = 30.39$, $MSE = 1,584.53$,
 1036 $p < .001$, $\hat{\eta}_G^2 = .159$; for the other-referential trials, the effect of valence is weaker,
 1037 $F(1.86, 63.08) = 2.85$, $MSE = 2,224.30$, $p = .069$, $\hat{\eta}_G^2 = .024$. We then focused on the self
 1038 conditions: the good-self condition (713 ± 12) is faster than neutral- (776 ± 11.8), $t(34) =$
 1039 -7.396 , $p < .0001$, and bad-self (772 ± 10.1) conditions, $t(34) = -5.66$, $p < .0001$. But
 1040 there is not difference between neutral- and bad-self conditions, $t(34) = 0.481$, $p = 0.881$.

1041 For the nonmatch trials, we didn't found any strong effect: identity, $F(1, 34) = 3.43$,
 1042 $MSE = 660.02$, $p = .073$, $\hat{\eta}_G^2 = .004$, valence $F(1.89, 64.33) = 0.40$, $MSE = 444.10$,
 1043 $p = .661$, $\hat{\eta}_G^2 = .001$, or interaction between the two $F(1.94, 66.02) = 2.42$, $MSE = 817.35$,
 1044 $p = .099$, $\hat{\eta}_G^2 = .007$.

1045 **BGLM.**

1046 *Signal detection theory analysis of accuracy.*

1047 We found that the d prime is greater when shapes were associated with good self
 1048 condition than with neutral self or bad self, but shapes associated with bad self and neutral
 1049 self didn't show differences. Comparing the self vs other under three condition revealed
 1050 that shapes associated with good self is greater than with good other, but with a weak
 1051 evidence. In contrast, for both neutral and bad valence condition, shapes associated with
 1052 other had greater d prime than with self.

1053 *Reaction time.*

1054 In reaction times, we found that same trends in the match trials as in the RT: while
 1055 the shapes associated with good self was greater than with good other (log mean diff =
 1056 -0.02858 , 95%HPD[-0.070898, 0.0154]), the direction is reversed for neutral and negative

1057 condition. see Figure 17

1058 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et
1059 al., 2020). We estimated separate drift rate (v), non-decision time (T_0), and boundary
1060 separation (a) for each condition. We found that the shapes tagged with good person has
1061 higher drift rate and higher boundary separation than shapes tagged with both neutral and
1062 bad person. Also, the shapes tagged with neutral person has a higher drift rate than
1063 shapes tagged with bad person, but not for the boundary separation. Finally, we found
1064 that shapes tagged with bad person had longer non-decision time (see figure 18)).

1065 **Experiment 3b**

1066 In study 3a, participants had to remember 6 pairs of association, which cause high
1067 cognitive load during the whole experiment. To eliminate the influence of cognitive load, we
1068 conducted study 3b, in which participant learn three aspect of self and stranger separately
1069 in to consecutive task. We hypothesize that we will replicate the pattern of study 3a, i.e.,
1070 the effect of moral valence only occurs for self-relevant conditions. ### Method

1071 **Participants.**

1072 Study 3b were finished in 2017, at that time we have calculated that the effect size
1073 (Cohen's d) of good-person (or good-self) vs. bad-person (or bad-other) was between $0.47 \sim$
1074 0.53, based on the data from Tsinghua community in study 1a, 1b, 2, 3a, 4a, and 4b. Based
1075 on this effect size, we estimated that 54 participants would allow we to detect the effect
1076 size of Cohen's $= 0.5$ with 95% power and alpha = 0.05, using G*power 3.192 (Faul,
1077 Erdfelder, Buchner, & Lang, 2009). Therefore, we planned to stop after we arrived this
1078 number. During the data collected at Wenzhou University, 61 participants (45 females; 19
1079 to 25 years of age, age = 20.42 ± 1.77) came to the testing room and we tested all of them
1080 during a single day. All participants were right-handed, and all had normal or
1081 corrected-to-normal vision. Informed consent was obtained from all participants prior to

1082 the experiment according to procedures approved by a local ethics committee. 4
1083 participants' data were excluded from analysis because their over all accuracy was lower
1084 than 60%, 1 more participant was excluded because of zero hit rate for one condition,
1085 leaving 56 participants (43 females; 19 to 25 years old, age = 20.27 ± 1.60).

1086 ***Design.***

1087 Study 3b has the same experimental design as 3a, with a $2 \times 3 \times 2$ within-subject
1088 design. The first variable was self-relevance, include two levels: self-relevant
1089 vs. stranger-relevant; the second variable was moral valence, include good, neutral and bad;
1090 the third variable was the matching between shape and label: match vs. mismatch.
1091 Stimuli. The stimuli used in study 3b share the same parameters with experiment 3a. 6
1092 shapes were included (triangle, square, circle, trapezoid, diamond, regular pentagon), as
1093 well as 6 labels, but the labels changed to "good self", "neutral self", "bad self", "good
1094 him/her", "bad him/her", "neutral him/her", the stranger's label is consistent with
1095 participants' gender. Same as study 3a, we asked participant to chosen an unfamiliar name
1096 of their own gender to be the stranger before showing them the relationship. Note, because
1097 of implementing error, the personal distance data did not collect for this experiment.

1098 ***Stimuli.***

1099 The stimuli used in study 3b is the same as in experiment 3a.

1100 ***Procedure.***

1101 In this experiment, participants finished two matching tasks, i.e., self-matching task,
1102 and other-matching task. In the self-matching task, participants first associate the three
1103 aspects of self to three different shapes, and then perform the matching task. In the
1104 other-matching task, participants first associate the three aspects of the stranger to three
1105 different shapes, and then perform the matching task. The order of self-task and other-task
1106 are counter-balanced among participants. Different from experiment 3a, after presenting
1107 the stimuli pair for 100ms, participant has 1900 ms to response, and they feedback with

₁₁₀₈ both accuracy and reaction time. As in study 3a, before each task, the instruction showed
₁₁₀₉ the meaning of each label to participants. The self-matching task and other-matching task
₁₁₁₀ were randomized between participants. Each participant finished 6 blocks, each have 120
₁₁₁₁ trials.

₁₁₁₂ ***Data Analysis.***

₁₁₁₃ Same as experiment 3a.

₁₁₁₄ **Results.**

₁₁₁₅ ***NHST.***

₁₁₁₆ Figure 19 shows d prime and reaction times of experiment 3b. Less than 5% correct
₁₁₁₇ trials with less than 200ms reaction times were excluded.

₁₁₁₈ *d prime.*

₁₁₁₉ There was no evidence for the main effect of valence, $F(1.92, 105.43) = 1.90$,
₁₁₂₀ $MSE = 0.33$, $p = .157$, $\hat{\eta}_G^2 = .005$, but we found a main effect of self-relevance,
₁₁₂₁ $F(1, 55) = 4.65$, $MSE = 0.89$, $p = .035$, $\hat{\eta}_G^2 = .017$, as well as the interaction,
₁₁₂₂ $F(1.90, 104.36) = 5.58$, $MSE = 0.26$, $p = .006$, $\hat{\eta}_G^2 = .011$.

₁₁₂₃ We then conducted separated ANOVA for self-referential and other-referential trials.
₁₁₂₄ The valence effect was shown for the self-referential conditions, $F(1.75, 96.42) = 6.73$,
₁₁₂₅ $MSE = 0.30$, $p = .003$, $\hat{\eta}_G^2 = .037$. Post-hoc test revealed that the Good-Self condition
₁₁₂₆ (2.15 ± 0.12) was with greater d prime than Neutral condition (1.83 ± 0.12 , $t(34) = 3.36$,
₁₁₂₇ $p = 0.0031$), and Bad-self condition (1.87 ± 0.12), $t(34) = 2.955$, $p = 0.01$. There was
₁₁₂₈ difference between neutral and bad condition, $t(34) = -0.039$, $p = 0.914$. However, no effect
₁₁₂₉ of valence was found for the other-referential condition $F(1.93, 105.97) = 0.61$,
₁₁₃₀ $MSE = 0.31$, $p = .539$, $\hat{\eta}_G^2 = .002$.

₁₁₃₁ *Reaction time.*

₁₁₃₂ We found interaction between Matchness and Valence ($F(1.86, 102.47) = 15.44$,

₁₁₃₃ $MSE = 3,112.78, p < .001, \hat{\eta}_G^2 = .006$) and then analyzed the matched trials and
₁₁₃₄ nonmatch trials separately, as in previous experiments.

₁₁₃₅ For the match trials, we found that the interaction between identity and valence,
₁₁₃₆ $F(1.67, 92.11) = 6.14, MSE = 6,472.48, p = .005, \hat{\eta}_G^2 = .009$, as well as the main effect of
₁₁₃₇ valence $F(1.88, 103.65) = 24.25, MSE = 5,994.25, p < .001, \hat{\eta}_G^2 = .038$, but not the effect
₁₁₃₈ of identity $F(1, 55) = 48.49, MSE = 25,892.59, p < .001, \hat{\eta}_G^2 = .153$. As for the d prime,
₁₁₃₉ we separated analyzed the self-referential and other-referential trials. For the
₁₁₄₀ Self-referential trials, we found the main effect of valence, $F(1.66, 91.38) = 23.98,$
₁₁₄₁ $MSE = 6,965.61, p < .001, \hat{\eta}_G^2 = .100$; for the other-referential trials, the effect of valence
₁₁₄₂ is weaker, $F(1.89, 103.94) = 5.96, MSE = 5,589.90, p = .004, \hat{\eta}_G^2 = .014$. We then focused
₁₁₄₃ on the self conditions: the good-self condition (713 ± 12) is faster than neutral- ($776 \pm$
₁₁₄₄ 11.8), $t(34) = -7.396, p < .0001$, and bad-self (772 ± 10.1) conditions, $t(34) = -5.66, p <$
₁₁₄₅ $.0001$. But there is not difference between neutral- and bad-self conditions, $t(34) = 0.481, p$
₁₁₄₆ $= 0.881$.

₁₁₄₇ For the nonmatch trials, we didn't found any strong effect: identity, $F(1, 55) = 10.31,$
₁₁₄₈ $MSE = 24,590.52, p = .002, \hat{\eta}_G^2 = .035$, valence $F(1.98, 108.63) = 20.57, MSE = 2,847.51,$
₁₁₄₉ $p < .001, \hat{\eta}_G^2 = .016$, or interaction between the two $F(1.93, 106.25) = 35.51,$
₁₁₅₀ $MSE = 1,939.88, p < .001, \hat{\eta}_G^2 = .019$.

₁₁₅₁ **BGLM.**

₁₁₅₂ *Signal detection theory analysis of accuracy.*

₁₁₅₃ We found that the d prime is greater when shapes were associated with good self
₁₁₅₄ condition than with neutral self or bad self, but shapes associated with bad self and neutral
₁₁₅₅ self didn't show differences. comparing the self vs other under three condition revealed that
₁₁₅₆ shapes associated with good self is greater than with good other, but with a weak evidence.
₁₁₅₇ In contrast, for both neutral and bad valence condition, shapes associated with other had
₁₁₅₈ greater d prime than with self.

1159 *Reaction time.*

1160 In reaction times, we found that same trends in the match trials as in the RT: while
1161 the shapes associated with good self was greater than with good other (log mean diff =
1162 -0.02858, 95%HPD[-0.070898, 0.0154]), the direction is reversed for neutral and negative
1163 condition. see Figure 20

1164 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et
1165 al., 2020). We estimated separate drift rate (v), non-decision time (T_0), and boundary
1166 separation (a) for each condition. We found that, similar to experiment 3a, the shapes
1167 tagged with good person has higher drift rate and higher boundary separation than shapes
1168 tagged with both neutral and bad person, but only for the self-referential condition. Also,
1169 the shapes tagged with neutral person has a higher drift rate than shapes tagged with bad
1170 person, but not for the boundary separation, and this effect also exist only for the
1171 self-referential condition.

1172 Interestingly, we found that in both self-referential and other-referential conditions,
1173 the shapes associated bad valence have higher drift rate and higher boundary separation.
1174 which might suggest that the shape associated with bad stimuli might be prioritized in the
1175 non-match trials (see figure 21)).

1176 **Experiment 6b**

1177 Experiment 6b was conducted to study the neural correlates of the prioritization
1178 effect of positive self, i.e., the neural underlying of the behavioral effect found int
1179 experiment 3a. However, as in experiment 6a, the procedure of this experiment was
1180 modified to adopted to ERP experiment.

1181 **Method.**

1182 **Participants.**

1183 23 college students (8 female, age = 22.86 ± 2.47) participated the current study, all
1184 of them were recruited from Tsinghua University in 2016. Informed consent was obtained
1185 from all participants prior to the experiment according to procedures approved by a local
1186 ethics committee. For day 1's data, 1 participant was excluded from the current analysis
1187 because of lower than 60% overall accuracy, remaining 22 participants (8 female, age =
1188 22.76 ± 2.49). For day 2's data, one participant dropped out, leaving 22 participants (9
1189 female, age = 23.05 ± 2.46), all of them has overall accuracy higher than 60%.

1190 ***Design.***

1191 The experimental design of this experiment is same as experiment 3: a $2 \times 3 \times 2$
1192 within-subject design with self-relevance (self-relevant vs. other-relevant), moral valence
1193 (good, neutral, and bad) and matchness between shape and label (match vs. mismatch) as
1194 within-subject variables.

1195 ***Stimuli.***

1196 As in experiment 3a, 6 shapes were included (triangle, square, circle, trapezoid,
1197 diamond, regular pentagon), as well as 6 labels (good self, neutral self, bad self, good
1198 person, bad person, neutral person). To match the concreteness of the label, we asked
1199 participant to chosen an unfamiliar name of their own gender to be the stranger.

1200 ***Procedure.***

1201 The procedure was similar to Experiment 2 and 6a. Subjects first learned the
1202 associations between labels and shapes and then completed a shape-label matching task. In
1203 each trial of the matching task, a fixation were first presented for 500 ms, followed by a 50
1204 ms label; then, a scrambled picture presented 900 ms. After the backward mask, the shape
1205 were presented on a noisy background for 50ms. Participant have to response in 1000ms
1206 after the presentation of the shape, and finally, a feedback screen was presented for 500 ms.
1207 The inter-trial interval (ITI) were randomly varied at the range of 1000 ~ 1400 ms.

1208 All the stimuli were presented on a gray background (RGB: 127, 127, 127). E-primed

1209 2.0 was used to present stimuli and collect behavioral results. Data were collected and
1210 analyzed when accuracy performance in total reached 60%.

1211 Because learning 6 associations was more difficult than 3 associations and participant
1212 might have low accuracy (see experiment 3a), the current study had extended to a two-day
1213 paradigm to maximizing the accurate trials that can be used in EEG data. At the first day,
1214 participants learnt the associations and finished 9 blocks of the matching task, each had
1215 120 trials, without EEG recording. That is, each condition has 90 trials.

1216 Participants came back to lab at the second day and finish the same task again, with
1217 EEG recorded. Before the EEG experiment, each participant finished a practice session
1218 again, if their accuracy is equal or higher than 85%, they start the experiment (one
1219 participant used lower threshold 75%). Each participant finished 18 blocks, each has 90
1220 trials. One participant finished additional 6 blocks because of high error rate at the
1221 beginning, another two participant finished addition 3 blocks because of the technique
1222 failure in recording the EEG data. To increase the number of trials that can be used for
1223 EEG data analysis, matched trials has twice number as mismatched trials, therefore, for
1224 matched trials each participants finished 180 trials for each condition, for mismatched
1225 trials, each conditions has 90 trials.

1226 ***Data Analysis.***

1227 Same as experiment 3a.

1228 **Results of Day 1.**

1229 ***NHST.***

1230 Figure 22 shows d' prime and reaction times of experiment 3b. Less than 5% correct
1231 trials with less than 200ms reaction times were excluded.

1232 ***d' prime.***

1233 There was no evidence for the main effect of valence, $F(1.91, 40.20) = 11.98,$

₁₂₃₄ $MSE = 0.15$, $p < .001$, $\hat{\eta}_G^2 = .040$, but we found a main effect of self-relevance,

₁₂₃₅ $F(1, 21) = 1.21$, $MSE = 0.20$, $p = .284$, $\hat{\eta}_G^2 = .003$, as well as the interaction,

₁₂₃₆ $F(1.28, 26.90) = 12.88$, $MSE = 0.21$, $p = .001$, $\hat{\eta}_G^2 = .041$.

₁₂₃₇ We then conducted separated ANOVA for self-referential and other-referential trials.

₁₂₃₈ The valence effect was shown for the self-referential conditions, $F(1.73, 36.42) = 29.31$,

₁₂₃₉ $MSE = 0.14$, $p < .001$, $\hat{\eta}_G^2 = .147$. Post-hoc test revealed that the Good-Self condition

₁₂₄₀ (2.15 ± 0.12) was with greater d prime than Neutral condition (1.83 ± 0.12 , $t(34) = 3.36$,

₁₂₄₁ $p = 0.0031$), and Bad-self condition (1.87 ± 0.12), $t(34) = 2.955$, $p = 0.01$. There was

₁₂₄₂ difference between neutral and bad condition, $t(34) = -0.039$, $p = 0.914$. However, no effect

₁₂₄₃ of valence was found for the other-referential condition $F(1.75, 36.72) = 0.00$, $MSE = 0.18$,

₁₂₄₄ $p = .999$, $\hat{\eta}_G^2 = .000$.

₁₂₄₅ *Reaction time.*

₁₂₄₆ We found interaction between Matchness and Valence ($F(1.79, 37.63) = 4.07$,

₁₂₄₇ $MSE = 704.90$, $p = .029$, $\hat{\eta}_G^2 = .003$) and then analyzed the matched trials and nonmatch

₁₂₄₈ trials separately, as in previous experiments.

₁₂₄₉ For the match trials, we found that the interaction between identity and valence,

₁₂₅₀ $F(1.72, 36.16) = 4.55$, $MSE = 1,560.90$, $p = .022$, $\hat{\eta}_G^2 = .015$, as well as the main effect of

₁₂₅₁ valence $F(1.93, 40.55) = 9.83$, $MSE = 1,951.84$, $p < .001$, $\hat{\eta}_G^2 = .044$, but not the effect of

₁₂₅₂ identity $F(1, 21) = 4.87$, $MSE = 2,032.05$, $p = .039$, $\hat{\eta}_G^2 = .012$. As for the d prime, we

₁₂₅₃ separated analyzed the self-referential and other-referential trials. For the Self-referential

₁₂₅₄ trials, we found the main effect of valence, $F(1.92, 40.38) = 14.48$, $MSE = 1,647.20$,

₁₂₅₅ $p < .001$, $\hat{\eta}_G^2 = .112$; for the other-referential trials, the effect of valence is weaker,

₁₂₅₆ $F(1.79, 37.50) = 1.04$, $MSE = 1,842.07$, $p = .356$, $\hat{\eta}_G^2 = .008$. We then focused on the self

₁₂₅₇ conditions: the good-self condition (713 ± 12) is faster than neutral- (776 ± 11.8), $t(34) =$

₁₂₅₈ -7.396 , $p < .0001$, and bad-self (772 ± 10.1) conditions, $t(34) = -5.66$, $p < .0001$. But

₁₂₅₉ there is not difference between neutral- and bad-self conditions, $t(34) = 0.481$, $p = 0.881$.

1260 For the nonmatch trials, we didn't found any strong effect: identity, $F(1, 21) = 2.76$,
 1261 $MSE = 1,718.93$, $p = .112$, $\hat{\eta}_G^2 = .006$, valence $F(1.61, 33.77) = 3.81$, $MSE = 1,532.21$,
 1262 $p = .041$, $\hat{\eta}_G^2 = .012$, or interaction between the two $F(1.90, 39.97) = 2.23$, $MSE = 720.80$,
 1263 $p = .123$, $\hat{\eta}_G^2 = .004$.

1264 **BGLM.**

1265 *Signal detection theory analysis of accuracy.*

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

1272 *Reaction time.*

1273 In reaction times, we found that same trends in the match trials as in the RT: while
 1274 the shapes associated with good self was greater than with good other ($\log \text{mean diff} =$
 1275 -0.02858 , $95\% \text{HPD}[-0.070898, 0.0154]$), the direction is reversed for neutral and negative
 1276 condition. see Figure 23

1277 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et
 1278 al., 2020). We estimated separate drift rate (v), non-decision time (T_0), and boundary
 1279 separation (a) for each condition. We found that, similar to experiment 3a, the shapes
 1280 tagged with good person has higher drift rate and higher boundary separation than shapes
 1281 tagged with both neutral and bad person, but only for the self-referential condition. Also,
 1282 the shapes tagged with neutral person has a higher drift rate than shapes tagged with bad
 1283 person, but not for the boundary separation, and this effect also exist only for the
 1284 self-referential condition.

1285 Interestingly, we found that in both self-referential and other-referential conditions,

1286 the shapes associated bad valence have higher drift rate and higher boundary separation.
1287 which might suggest that the shape associated with bad stimuli might be prioritized in the
1288 non-match trials (see figure 24).

1289 **Part 3: Implicit binding between valence and identity**

1290 In this part, we reported two studies in which the moral valence or the self-referential
1291 processing is not task-relevant. We are interested in testing whether the task-relevance will
1292 eliminate the effect observed in previous experiment.

1293 **Experiment 4a: Morality as task-irrelevant variable**

1294 In part two (experiment 3a and 3b), participants learned the association between self
1295 and moral valence directly. In Experiment 4a, we examined whether the interaction
1296 between moral valence and identity occur even when one of the variable was irrelevant to
1297 the task. In experiment 4a, participants learnt associations between shapes and self/other
1298 labels, then made perceptual match judgments only about the self or other conditions
1299 labels and shapes (cf. Sui et al. (2012)). However, we presented labels of different moral
1300 valence in the shapes, which means that the moral valence factor become task irrelevant. If
1301 the binding between moral good and self is intrinsic and automatic, then we will observe
1302 that facilitating effect of moral good for self conditions, but not for other conditions.

1303 **Method.**

1304 ***Participants.***

1305 64 participants (37 female, age = 19.70 ± 1.22) participated the current study, 32 of
1306 them were from Tsinghua University in 2015, 32 were from Wenzhou University
1307 participated in 2017. All participants were right-handed, and all had normal or
1308 corrected-to-normal vision. Informed consent was obtained from all participants prior to
1309 the experiment according to procedures approved by a local ethics committee. The data

1310 from 5 participants from Wenzhou site were excluded from analysis because their accuracy
1311 was close to chance (< 0.6). The results for the remaining 59 participants (33 female, age
1312 = 19.78 ± 1.20) were analyzed and reported.

1313 ***Design.***

1314 As in Experiment 3, a $2 \times 3 \times 2$ within-subject design was used. The first variable was
1315 self-relevance (self and stranger associations); the second variable was moral valence (good,
1316 neutral and bad associations); the third variable was the matching between shape and label
1317 (matching vs. non-match for the personal association). However, in this the task,
1318 participants only learn the association between two geometric shapes and two labels (self
1319 and other), i.e., only self-relevance were related to the task. The moral valence
1320 manipulation was achieved by embedding the personal label of the labels in the geometric
1321 shapes, see below. For simplicity, the trials where shapes where paired with self and with a
1322 word of “good person” inside were shorted as good-self condition, similarly, the trials where
1323 shapes paired with the self and with a word of “bad person” inside were shorted as bad-self
1324 condition. Hence, we also have six conditions: good-self, neutral-self, bad-self, good-other,
1325 neutral-other, and bad-other.

1326 ***Stimuli.***

1327 2 shapes were included (circle, square) and each appeared above a central fixation
1328 cross with the personal label appearing below. However, the shapes were not empty but
1329 with a two-Chinese-character word in the middle, the word was one of three labels with
1330 different moral valence: “good person”, “bad person” and “neutral person”. Before the
1331 experiment, participants learned the self/other association, and were informed to only
1332 response to the association between shapes’ configures and the labels below the fixation, but
1333 ignore the words within shapes. Besides the behavioral experiments, participants from
1334 Tsinghua community also finished questionnaires as Experiments 3, and participants from
1335 Wenzhou community finished a series of questionnaire as the other experiment finished in

₁₃₃₆ Wenzhou.

₁₃₃₇ ***Procedure.***

₁₃₃₈ The procedure was similar to Experiment 1. There were 6 blocks of trial, each with
₁₃₃₉ 120 trials for 2017 data. Due to procedure error, the data collected in 2015 in Tsinghua
₁₃₄₀ community only have 60 trials for each block, i.e., 30 trials per condition.

₁₃₄₁ As in study 3a, before each task, the instruction showed the meaning of each label to
₁₃₄₂ participants. The self-matching task and other-matching task were randomized between
₁₃₄₃ participants. Each participant finished 6 blocks, each have 120 trials.

₁₃₄₄ ***Data Analysis.***

₁₃₄₅ Same as experiment 3a.

₁₃₄₆ **Results.**

₁₃₄₇ ***NHST.***

₁₃₄₈ Figure 25 shows d prime and reaction times of experiment 3a. Less than 5% correct
₁₃₄₉ trials with less than 200ms reaction times were excluded.

₁₃₅₀ d prime.

₁₃₅₁ There was no evidence for the main effect of valence, $F(1.93, 111.66) = 0.53$,
₁₃₅₂ $MSE = 0.12$, $p = .581$, $\hat{\eta}_G^2 = .000$, but we found a main effect of self-relevance,
₁₃₅₃ $F(1, 58) = 121.04$, $MSE = 0.48$, $p < .001$, $\hat{\eta}_G^2 = .189$, as well as the interaction,
₁₃₅₄ $F(1.99, 115.20) = 4.12$, $MSE = 0.14$, $p = .019$, $\hat{\eta}_G^2 = .004$.

₁₃₅₅ We then conducted separated ANOVA for self-referential and other-referential trials.
₁₃₅₆ The valence effect was shown for the self-referential conditions, $F(1.95, 112.92) = 3.01$,
₁₃₅₇ $MSE = 0.15$, $p = .055$, $\hat{\eta}_G^2 = .008$. Post-hoc test revealed that the Good-Self condition
₁₃₅₈ (2.15 ± 0.12) was with greater d prime than Neutral condition (1.83 ± 0.12 , $t(34) = 3.36$,
₁₃₅₉ $p = 0.0031$), and Bad-self condition (1.87 ± 0.12), $t(34) = 2.955$, $p = 0.01$. There was
₁₃₆₀ difference between neutral and bad condition, $t(34) = -0.039$, $p = 0.914$. However, no effect

₁₃₆₁ of valence was found for the other-referential condition $F(1.98, 114.61) = 1.75$,

₁₃₆₂ $MSE = 0.10$, $p = .179$, $\hat{\eta}_G^2 = .003$.

₁₃₆₃ *Reaction time.*

₁₃₆₄ We found interaction between Matchness and Valence ($F(1.94, 112.64) = 0.84$,

₁₃₆₅ $MSE = 465.35$, $p = .432$, $\hat{\eta}_G^2 = .000$) and then analyzed the matched trials and nonmatch

₁₃₆₆ trials separately, as in previous experiments.

₁₃₆₇ For the match trials, we found that the interaction between identity and valence,

₁₃₆₈ $F(1.90, 110.18) = 4.41$, $MSE = 465.91$, $p = .016$, $\hat{\eta}_G^2 = .003$, as well as the main effect of

₁₃₆₉ valence $F(1.98, 114.82) = 0.94$, $MSE = 606.30$, $p = .392$, $\hat{\eta}_G^2 = .001$, but not the effect of

₁₃₇₀ identity $F(1, 58) = 124.15$, $MSE = 4,037.53$, $p < .001$, $\hat{\eta}_G^2 = .257$. As for the d prime, we

₁₃₇₁ separated analyzed the self-referential and other-referential trials. For the Self-referential

₁₃₇₂ trials, we found the main effect of valence, $F(1.97, 114.32) = 6.29$, $MSE = 367.25$,

₁₃₇₃ $p = .003$, $\hat{\eta}_G^2 = .006$; for the other-referential trials, the effect of valence is weaker,

₁₃₇₄ $F(1.95, 112.89) = 0.35$, $MSE = 699.50$, $p = .699$, $\hat{\eta}_G^2 = .001$. We then focused on the self

₁₃₇₅ conditions: the good-self condition (713 ± 12) is faster than neutral- (776 ± 11.8), $t(34) =$

₁₃₇₆ -7.396 , $p < .0001$, and bad-self (772 ± 10.1) conditions, $t(34) = -5.66$, $p < .0001$. But

₁₃₇₇ there is not difference between neutral- and bad-self conditions, $t(34) = 0.481$, $p = 0.881$.

₁₃₇₈ For the nonmatch trials, we didn't found any strong effect: identity, $F(1, 58) = 0.16$,

₁₃₇₉ $MSE = 1,547.37$, $p = .692$, $\hat{\eta}_G^2 = .000$, valence $F(1.96, 113.52) = 0.68$, $MSE = 390.26$,

₁₃₈₀ $p = .508$, $\hat{\eta}_G^2 = .000$, or interaction between the two $F(1.90, 110.27) = 0.04$,

₁₃₈₁ $MSE = 585.80$, $p = .953$, $\hat{\eta}_G^2 = .000$.

₁₃₈₂ **BGLM.**

₁₃₈₃ *Signal detection theory analysis of accuracy.*

₁₃₈₄ We found that the d prime is greater when shapes were associated with good self

₁₃₈₅ condition than with neutral self or bad self, but shapes associated with bad self and neutral

₁₃₈₆ self didn't show differences. comparing the self vs other under three condition revealed that

1387 shapes associated with good self is greater than with good other, but with a weak evidence.
1388 In contrast, for both neutral and bad valence condition, shapes associated with other had
1389 greater d' prime than with self.

1390 *Reaction time.*

1391 In reaction times, we found that same trends in the match trials as in the RT: while
1392 the shapes associated with good self was greater than with good other (log mean diff =
1393 -0.02858, 95%HPD[-0.070898, 0.0154]), the direction is reversed for neutral and negative
1394 condition. see Figure 26

1395 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et
1396 al., 2020). We estimated separate drift rate (v), non-decision time (T_0), and boundary
1397 separation (a) for each condition. We found that the shapes tagged with good person has
1398 higher drift rate and higher boundary separation than shapes tagged with both neutral and
1399 bad person. Also, the shapes tagged with neutral person has a higher drift rate than
1400 shapes tagged with bad person, but not for the boundary separation. Finally, we found
1401 that shapes tagged with bad person had longer non-decision time (see figure 27)).

1402 **Experiment 4b: Morality as task-irrelevant variable**

1403 In study 4b, we changed the role of valence and identity in task. In this experiment,
1404 participants learn the association between moral valence and the made perceptual match
1405 judgments to associations between different moral valence and shapes as in study 1-3.
1406 Different from experiment 1 ~ 3, we made put the labels of “self/other” in the shapes so
1407 that identity served as an task irrelevant variable. As in experiment 4b, we also
1408 hypothesized that the intrinsic binding between morally good and self will enhance the
1409 performance of good self condition, even identity is irrelevant to the task.

1410 **Method.**

Participants.

53 participants (39 female, age = 20.57 ± 1.81) participated the current study, 34 of them were from Tsinghua University in 2015, 19 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 8 participants from Wenzhou site were excluded from analysis because their accuracy was close to chance (< 0.6). The results for the remaining 45 participants (33 female, age = 20.78 ± 1.76) 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 the 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

¹⁴³⁷ different moral valence: “good person”, “bad person” and “neutral person”. Before the
¹⁴³⁸ experiment, participants learned the self/other association, and were informed to only
¹⁴³⁹ response to the association between shapes’ configures and the labels below the fixation, but
¹⁴⁴⁰ ignore the words within shapes. Besides the behavioral experiments, participants from
¹⁴⁴¹ Tsinghua community also finished questionnaires as Experiments 3, and participants from
¹⁴⁴² Wenzhou community finished a series of questionnaire as the other experiment finished in
¹⁴⁴³ Wenzhou.

¹⁴⁴⁴ ***Procedure.***

¹⁴⁴⁵ The procedure was similar to Experiment 1. There were 6 blocks of trial, each with
¹⁴⁴⁶ 120 trials for 2017 data. Due to procedure error, the data collected in 2015 in Tsinghua
¹⁴⁴⁷ community only have 60 trials for each block, i.e., 30 trials per condition.

¹⁴⁴⁸ As in study 3a, before each task, the instruction showed the meaning of each label to
¹⁴⁴⁹ participants. The self-matching task and other-matching task were randomized between
¹⁴⁵⁰ participants. Each participant finished 6 blocks, each have 120 trials.

¹⁴⁵¹ ***Data Analysis.***

¹⁴⁵² Same as experiment 3a.

¹⁴⁵³ ***Results.***

¹⁴⁵⁴ ***NHST.***

¹⁴⁵⁵ Figure 28 shows d prime and reaction times of experiment 3a. Less than 5% correct
¹⁴⁵⁶ trials with less than 200ms reaction times were excluded.

¹⁴⁵⁷ d prime.

¹⁴⁵⁸ There was no evidence for the main effect of valence, $F(1.59, 69.94) = 2.34$,
¹⁴⁵⁹ $MSE = 0.48$, $p = .115$, $\hat{\eta}_G^2 = .010$, but we found a main effect of self-relevance,
¹⁴⁶⁰ $F(1, 44) = 0.00$, $MSE = 0.08$, $p = .994$, $\hat{\eta}_G^2 = .000$, as well as the interaction,
¹⁴⁶¹ $F(1.96, 86.41) = 0.53$, $MSE = 0.10$, $p = .585$, $\hat{\eta}_G^2 = .001$.

1462 We then conducted separated ANOVA for self-referential and other-referential trials.

1463 The valence effect was shown for the self-referential conditions, $F(1.75, 76.86) = 3.08$,

1464 $MSE = 0.25$, $p = .058$, $\hat{\eta}_G^2 = .017$. Post-hoc test revealed that the Good-Self condition

1465 (2.15 ± 0.12) was with greater d prime than Neutral condition (1.83 ± 0.12 , $t(34) = 3.36$,

1466 $p = 0.0031$), and Bad-self condition (1.87 ± 0.12), $t(34) = 2.955$, $p = 0.01$. There was

1467 difference between neutral and bad condition, $t(34) = -0.039$, $p = 0.914$. However, no effect

1468 of valence was found for the other-referential condition $F(1.63, 71.50) = 1.07$, $MSE = 0.33$,

1469 $p = .336$, $\hat{\eta}_G^2 = .006$.

1470 *Reaction time.*

1471 We found interaction between Matchness and Valence ($F(1.87, 82.50) = 18.58$,

1472 $MSE = 1,291.12$, $p < .001$, $\hat{\eta}_G^2 = .023$) and then analyzed the matched trials and

1473 nonmatch trials separately, as in previous experiments.

1474 For the match trials, we found that the interaction between identity and valence,

1475 $F(1.86, 81.84) = 5.22$, $MSE = 308.30$, $p = .009$, $\hat{\eta}_G^2 = .003$, as well as the main effect of

1476 valence $F(1.80, 79.37) = 11.04$, $MSE = 2,937.54$, $p < .001$, $\hat{\eta}_G^2 = .059$, but not the effect of

1477 identity $F(1, 44) = 0.23$, $MSE = 263.26$, $p = .632$, $\hat{\eta}_G^2 = .000$. As for the d prime, we

1478 separated analyzed the self-referential and other-referential trials. For the Self-referential

1479 trials, we found the main effect of valence, $F(1.74, 76.48) = 13.69$, $MSE = 1,732.08$,

1480 $p < .001$, $\hat{\eta}_G^2 = .079$; for the other-referential trials, the effect of valence is weaker,

1481 $F(1.87, 82.44) = 7.09$, $MSE = 1,527.43$, $p = .002$, $\hat{\eta}_G^2 = .043$. We then focused on the self

1482 conditions: the good-self condition (713 ± 12) is faster than neutral- (776 ± 11.8), $t(34) =$

1483 -7.396 , $p < .0001$, and bad-self (772 ± 10.1) conditions, $t(34) = -5.66$, $p < .0001$. But

1484 there is not difference between neutral- and bad-self conditions, $t(34) = 0.481$, $p = 0.881$.

1485 For the nonmatch trials, we didn't found any strong effect: identity, $F(1, 44) = 1.96$,

1486 $MSE = 319.47$, $p = .169$, $\hat{\eta}_G^2 = .001$, valence $F(1.69, 74.54) = 6.59$, $MSE = 886.19$,

1487 $p = .004$, $\hat{\eta}_G^2 = .010$, or interaction between the two $F(1.88, 82.57) = 0.31$, $MSE = 316.96$,

¹⁴⁸⁸ $p = .718$, $\hat{\eta}_G^2 = .000$.

¹⁴⁸⁹ **BGLM.**

¹⁴⁹⁰ *Signal detection theory analysis of accuracy.*

¹⁴⁹¹ We found that the d prime is greater when shapes were associated with good self
¹⁴⁹² condition than with neutral self or bad self, but shapes associated with bad self and neutral
¹⁴⁹³ self didn't show differences. comparing the self vs other under three condition revealed that
¹⁴⁹⁴ shapes associated with good self is greater than with good other, but with a weak evidence.
¹⁴⁹⁵ In contrast, for both neutral and bad valence condition, shapes associated with other had
¹⁴⁹⁶ greater d prime than with self.

¹⁴⁹⁷ *Reaction time.*

¹⁴⁹⁸ In reaction times, we found that same trends in the match trials as in the RT: while
¹⁴⁹⁹ the shapes associated with good self was greater than with good other ($\log \text{mean diff} =$
¹⁵⁰⁰ -0.02858 , $95\% \text{HPD}[-0.070898, 0.0154]$), the direction is reversed for neutral and negative
¹⁵⁰¹ condition. see Figure 29

¹⁵⁰² **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et
¹⁵⁰³ al., 2020). We estimated separate drift rate (v), non-decision time (T_0), and boundary
¹⁵⁰⁴ separation (a) for each condition. We found that the shapes tagged with good person has
¹⁵⁰⁵ higher drift rate and higher boundary separation than shapes tagged with both neutral and
¹⁵⁰⁶ bad person. Also, the shapes tagged with neutral person has a higher drift rate than
¹⁵⁰⁷ shapes tagged with bad person, but not for the boundary separation. Finally, we found
¹⁵⁰⁸ that shapes tagged with bad person had longer non-decision time (see figure 30)).

1509

Results

1510 Effect of moral valence

1511 In this part, we synthesized results from experiment 1a, 1b, 1c, 2, 5 and 6a. Data
1512 from 192 participants were included in these analyses. We found differences between
1513 positive and negative conditions on RT was Cohen's $d = -0.58 \pm 0.06$, 95% CI [-0.70 -0.47];
1514 on d' was Cohen's $d = 0.24 \pm 0.05$, 95% CI [0.15 0.34]. The effect was also observed
1515 between positive and neutral condition, RT: Cohen's $d = -0.44 \pm 0.10$, 95% CI [-0.63
1516 -0.25]; d' : Cohen's $d = 0.31 \pm 0.07$, 95% CI [0.16 0.45]. And the difference between neutral
1517 and bad conditions are not significant, RT: Cohen's $d = 0.15 \pm 0.07$, 95% CI [0.00 0.30];
1518 d' : Cohen's $d = 0.07 \pm 0.07$, 95% CI [-0.08 0.21]. See Figure 31 left panel.

1519 Interaction between valence and self-reference

1520 In this part, we combined the experiments that explicitly manipulated the
1521 self-reference and valence, which includes 3a, 3b, 6b, 7a, and 7b. For the positive versus
1522 negative contrast, data were from five experiments with 178 participants; for positive
1523 versus neutral and neutral versus negative contrasts, data were from three experiments ((

1524 3a, 3b, and 6b) with 108 participants.

1525 In most of these experiments, the interaction between self-reference and valence was
1526 significant (see results of each experiment in supplementary materials). In the
1527 mini-meta-analysis, we analyzed the valence effect for self-referential condition and
1528 other-referential condition separately.

1529 For the self-referential condition, we found the same pattern as in the first part of
1530 results. That is we found significant differences between positive and neutral as well as
1531 positive and negative, but not neutral and negative. The effect size of RT between positive
1532 and negative is Cohen's $d = -0.89 \pm 0.12$, 95% CI [-1.11 -0.66]; on d' was Cohen's $d = 0.61$

1533 ± 0.09 , 95% CI [0.44 0.78]. The effect was also observed between positive and neutral
1534 condition, RT: Cohen's $d = -0.76 \pm 0.13$, 95% CI [-1.01 -0.50]; d' : Cohen's $d = 0.69 \pm$
1535 0.14, 95% CI [0.42 0.96]. And the difference between neutral and bad conditions are not
1536 significant, RT: Cohen's $d = 0.03 \pm 0.13$, 95% CI [-0.22 0.29]; d' : Cohen's $d = 0.08 \pm 0.08$,
1537 95% CI [-0.07 0.24]. See Figure 31 the middle panel.

1538 For the other-referential condition, we found that only the difference between positive
1539 and negative on RT was significant, all the other conditions were not. The effect size of RT
1540 between positive and negative is Cohen's $d = -0.28 \pm 0.05$, 95% CI [-0.38 -0.17]; on d' was
1541 Cohen's $d = -0.02 \pm 0.08$, 95% CI [-0.17 0.13]. The effect was not observed between
1542 positive and neutral condition, RT: Cohen's $d = -0.12 \pm 0.10$, 95% CI [-0.31 0.06]; d' :
1543 Cohen's $d = 0.01 \pm 0.08$, 95% CI [-0.16 0.17]. And the difference between neutral and bad
1544 conditions are not significant, RT: Cohen's $d = 0.14 \pm 0.09$, 95% CI [-0.03 0.31]; d' :
1545 Cohen's $d = 0.05 \pm 0.07$, 95% CI [-0.08 0.18]. See Figure 31 right panel.

1546 Generalizability of the valence effect

1547 In this part, we reported the results from experiment 4 in which either moral valence
1548 or self-reference were manipulated as task-irrelevant stimuli.

1549 For experiment 4a, when self-reference was the target and moral valence was
1550 task-irrelevant, we found that only under the implicit self-referential condition, i.e., when
1551 the moral words were presented as task irrelevant stimuli, there was the main effect of
1552 valence and interaction between valence and reference for both d prime and RT (See
1553 supplementary results for the detailed statistics). For d prime, we found good-self
1554 condition (2.55 ± 0.86) had higher d prime than bad-self condition (2.38 ± 0.80); good self
1555 condition was also higher than neutral self (2.45 ± 0.78) but there was not statistically
1556 significant, while the neutral-self condition was higher than bad self condition and not
1557 significant neither. For reaction times, good-self condition (654.26 ± 67.09) were faster

1558 relative to bad-self condition (665.64 ± 64.59), and over neutral-self condition ($664.26 \pm$
1559 64.71). The difference between neutral-self and bad-self conditions were not significant.
1560 However, for the other-referential condition, there was no significant differences between
1561 different valence conditions. See Figure 32.

1562 For experiment 4b, when valence was the target and the identity was task-irrelevant,
1563 we found a strong valence effect (see supplementary results and Figure 33, Figure 34).

1564 In this experiment, the advantage of good-self condition can only be disentangled by
1565 comparing the self-referential and other-referential conditions. Therefore, we calculated the
1566 differences between the valence effect under self-referential and other referential conditions
1567 and used the weighted variance as the variance of this differences. We found this
1568 modulation effect on RT. The valence effect of RT was stronger in self-referential than
1569 other-referential for the Good vs. Neutral condition (-0.33 ± 0.01), and to a less extent the
1570 Good vs. Bad condition (-0.17 ± 0.01). While the size of the other effect's CI included
1571 zero, suggestion those effects didn't differ from zero. See Figure 35.

1572 Specificity of valence effect

1573 In this part, we analyzed the results from experiment 5, which included positive,
1574 neutral, and negative valence from four different domains: morality, emotion, aesthetics of
1575 human, and aesthetics of scene. We found interaction between valence and domain for both
1576 *d* prime and RT (match trials). A common pattern appeared in all four domains: each
1577 domain showed a binary results instead of gradient on both *d* prime and RT. For morality,
1578 aesthetics of human, and aesthetics of scene, there was a positivity effect where the positive
1579 conditions had advantages over both neutral (greater *d* prime and faster RT), while neutral
1580 and negative conditions didn't differ from each other. But for the emotional stimuli, there
1581 was a reversed negativity effect: positive and neutral conditions were not significantly
1582 different from each other but both had advantage over negative conditions. See

1583 supplementary materials for detailed statistics. Also note that the effect size in moral
1584 domain is smaller than the aesthetic domains (beauty of people and beauty of scene). See
1585 Figure 36.

1586 **Self-reported personal distance**

1587 See Figure 37.

1588 **Correlation analyses**

1589 The reliability of questionnaires can be found in (Liu et al., 2020). We calculated the
1590 correlation between the data from behavioral task and the questionnaire data. First, we
1591 calculated the score for each scale based on their structure and factor loading, instead of
1592 sum score (McNeish & Wolf, 2020). Then, we used SEM to estimate the correlation
1593 because it can include measurement model and statistical model in a unified framework.

1594 To make sure that what we found were not false positive, we used two method to
1595 ensure the robustness of our analysis. first, we split the data into two half: the data with
1596 self and without, then, we used the conditional random forest to find the robust correlation
1597 in the exploratory data (with self reference) that can be replicated in the confirmatory data
1598 (without the self reference). The robust correlation were then analyzed using SEM

1599 Instead of use the exploratory correlation analysis, we used a more principled way to
1600 explore the correlation between parameter of HDDM (v , t , and a) and scale scores and
1601 person distance.

1602 We didn't find the correlation between scale scores and the parameters of HDDM,
1603 but found weak correlation between personal distance and the parameter estimated from
1604 Good and neutral conditions.

1605 First, boundary separation (a) of moral good condition was correlated with both
1606 Self-Bad distance ($r = 0.198$, 95% CI [], $p = 0.0063$) and Neutral-Bad distance

1607 ($r = 0.1571$, 95% CI [], $p = 0.031$). At the same time, the non-decision time is negatively
1608 correlated with Self-Bad distance ($r = 0.169$, 95% CI [], $p = 0.0197$). See Figure 38.

1609 Second, we found the boundary separation of neutral condition is positively
1610 correlated with the personal distance between self and good distance ($r = 0.189$, 95% CI [],
1611 $p = 0.036$), but negatively correlated with self-neutral distance($r = -0.183$, 95% CI [],
1612 $p = 0.042$). Also, the drift rate of the neutral condition is positively correlated with the
1613 Self-Bad distance ($r = 0.177$, 95% CI [], $p = 0.048$).a. See figure 39

1614 We also explored the correlation between behavioral data and questionnaire scores
1615 separately for experiments with and without self-referential, however, the sample size is
1616 very low for some conditions.

1617 Discussion

1618 References

- 1619 Abele, A. E., Ellemers, N., Fiske, S. T., Koch, A., & Yzerbyt, V. (2020). Navigating the
1620 social world: Toward an integrated framework for evaluating self, individuals, and
1621 groups. *Psychological Review*. <https://doi.org/10.1037/rev0000262>
- 1622 Amodio, D. M. (2019). Social cognition 2.0: An interactive memory systems account.
1623 *Trends in Cognitive Sciences*, 23(1), 21–33.
1624 <https://doi.org/10.1016/j.tics.2018.10.002>
- 1625 Anderson, E., Siegel, E. H., Bliss-Moreau, E., & Barrett, L. F. (2011). The visual impact
1626 of gossip. *Science*, 332(6036), 1446–1448. <https://doi.org/10.1126/science.1201574>
- 1627 Brainard, D. H. (1997). The psychophysics toolbox. *Spatial Vision*, 10(4), 433–436.
1628 Journal Article.
- 1629 Bruner, J. S. (1957). On perceptual readiness. *Psychological Review*, 64(2), 123–152.
1630 <https://doi.org/10.1037/h0043805>

- 1631 Bürkner, P.-C. (2017). Brms: An r package for bayesian multilevel models using stan.
1632 *Journal of Statistical Software; Vol 1, Issue 1 (2017)*. Journal Article. Retrieved
1633 from
1634 <https://www.jstatsoft.org/v080/i01>
1635 Ahttp://dx.doi.org/10.18637/jss.v080.i01
1636 Carlson, R. W., Maréchal, M. A., Oud, B., Fehr, E., & Crockett, M. J. (2020). Motivated
1637 misremembering of selfish decisions. *Nature Communications*, 11(1), 2100.
1638 <https://doi.org/10.1038/s41467-020-15602-4>
1639 Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M., ...
1640 Riddell, A. (2017). Stan: A probabilistic programming language. *Journal of
1641 Statistical Software*, 76(1). Journal Article. <https://doi.org/10.18637/jss.v076.i01>
1642 Cooper, H., Hedges, L. V., & Valentine, J. C. (2009). *The handbook of research synthesis
1643 and meta-analysis* (2nd ed.). Book, New York: Sage.
1644 DeCarlo, L. T. (1998). Signal detection theory and generalized linear models. *Psychological
1645 Methods*, 3(2), 186–205. Journal Article. <https://doi.org/10.1037/1082-989X.3.2.186>
1646 Dzhelyova, M., Perrett, D. I., & Jentzsch, I. (2012). Temporal dynamics of trustworthiness
1647 perception. *Brain Research*, 1435, 81–90.
1648 <https://doi.org/10.1016/j.brainres.2011.11.043>
1649 Enock, F., Sui, J., Hewstone, M., & Humphreys, G. W. (2018). Self and team prioritisation
1650 effects in perceptual matching: Evidence for a shared representation. *Acta
1651 Psychologica*, 182, 107–118. <https://doi.org/10.1016/j.actpsy.2017.11.011>
1652 Faul, F., Erdfelder, E., Buchner, A., & Lang, A.-G. (2009). Statistical power analyses using
1653 g*Power 3.1: Tests for correlation and regression analyses. *Behavior Research
1654 Methods*, 41(4), 1149–1160. <https://doi.org/10.3758/BRM.41.4.1149>
1655 Freeman, J. B., & Ambady, N. (2011). A dynamic interactive theory of person construal.
1656 *Psychological Review*, 118(2), 247–279. <https://doi.org/10.1037/a0022327>

- 1656 Gantman, A. P., & Van Bavel, J. J. (2014). The moral pop-out effect: Enhanced
1657 perceptual awareness of morally relevant stimuli. *Cognition*, 132(1), 22–29.
1658 <https://doi.org/10.1016/j.cognition.2014.02.007>
- 1659 Goh, J. X., Hall, J. A., & Rosenthal, R. (2016). Mini meta-analysis of your own studies:
1660 Some arguments on why and a primer on how. *Social and Personality Psychology
Compass*, 10(10), 535–549. Journal Article. <https://doi.org/10.1111/spc3.12267>
- 1662 Goodwin, G. P. (2015). Moral character in person perception. *Current Directions in
Psychological Science*, 24(1), 38–44. <https://doi.org/10.1177/0963721414550709>
- 1664 Goodwin, G. P., Piazza, J., & Rozin, P. (2014). Moral character predominates in person
1665 perception and evaluation. *Journal of Personality and Social Psychology*, 106(1),
1666 148–168. <https://doi.org/10.1037/a0034726>
- 1667 Hu, C.-P., Lan, Y., Macrae, C. N., & Sui, J. (2020). Good me bad me: Does valence
1668 influence self-prioritization during perceptual decision-making? *Collabra:
Psychology*, 6(1), 20. Journal Article. <https://doi.org/10.1525/collabra.301>
- 1670 Juechems, K., & Summerfield, C. (2019). Where does value come from? *Trends in
Cognitive Sciences*, 23(10), 836–850. <https://doi.org/10.1016/j.tics.2019.07.012>
- 1672 Kirby, K. N., & Gerlanc, D. (2013). BootES: An r package for bootstrap confidence
1673 intervals on effect sizes. *Behavior Research Methods*, 45(4), 905–927.
1674 <https://doi.org/10.3758/s13428-013-0330-5>
- 1675 Liu, Q., Wang, F., Yan, W., Peng, K., Sui, J., & Hu, C.-P. (2020). Questionnaire data from
1676 the revision of a chinese version of free will and determinism plus scale. *Journal of
Open Psychology Data*, 8(1), 1. Journal Article. <https://doi.org/10.5334/jopd.49/>
- 1678 Matzke, D., & Wagenmakers, E.-J. (2009). Psychological interpretation of the ex-gaussian
1679 and shifted wald parameters: A diffusion model analysis. *Psychonomic Bulletin &
Review*, 16(5), 798–817. <https://doi.org/10.3758/PBR.16.5.798>

- 1681 McNeish, D., & Wolf, M. G. (2020). Thinking twice about sum scores. *Behavior Research
1682 Methods*. <https://doi.org/10.3758/s13428-020-01398-0>
- 1683 Pelli, D. G. (1997). The videotoolbox software for visual psychophysics: Transforming
1684 numbers into movies. *Spatial Vision*, 10(4), 437–442. Journal Article.
- 1685 Reicher, S., & Hopkins, N. (2016). Perception, action, and the social dynamics of the
1686 variable self. *Psychological Inquiry*, 27(4), 341–347.
1687 <https://doi.org/10.1080/1047840X.2016.1217584>
- 1688 Rouder, J. N., & Lu, J. (2005). An introduction to bayesian hierarchical models with an
1689 application in the theory of signal detection. *Psychonomic Bulletin & Review*,
1690 12(4), 573–604. Journal Article. <https://doi.org/10.3758/bf03196750>
- 1691 Rousselet, G. A., & Wilcox, R. R. (2019). Reaction times and other skewed distributions:
1692 Problems with the mean and the median. *Meta-Psychology*. preprint.
1693 <https://doi.org/10.1101/383935>
- 1694 Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2013). Life after p-hacking. Conference
1695 Proceedings. <https://doi.org/10.2139/ssrn.2205186>
- 1696 Stanislaw, H., & Todorov, N. (1999). Calculation of signal detection theory measures.
1697 *Behavior Research Methods, Instruments, & Computers*, 31(1), 137–149. Journal
1698 Article. <https://doi.org/10.3758/BF03207704>
- 1699 Stanley, M. L., & De Brigard, F. (2019). Moral memories and the belief in the good self.
1700 *Current Directions in Psychological Science*, 28(4), 387–391.
1701 <https://doi.org/10.1177/0963721419847990>
- 1702 Stolier, R. M., & Freeman, J. B. (2016). Functional and temporal considerations for
1703 top-down influences in social perception. *Psychological Inquiry*, 27(4), 352–357.
1704 <https://doi.org/10.1080/1047840X.2016.1216034>
- 1705 Sui, J., He, X., & Humphreys, G. W. (2012). Perceptual effects of social salience: Evidence

- 1706 from self-prioritization effects on perceptual matching. *Journal of Experimental*
1707 *Psychology: Human Perception and Performance*, 38(5), 1105–1117. Journal
1708 Article. <https://doi.org/10.1037/a0029792>
- 1709 Tappin, B. M., & McKay, R. T. (2017). The illusion of moral superiority. *Social*
1710 *Psychological and Personality Science*, 8(6), 623–631.
1711 <https://doi.org/10.1177/1948550616673878>
- 1712 Verosky, S. C., & Todorov, A. (2010). Differential neural responses to faces physically
1713 similar to the self as a function of their valence. *NeuroImage*, 49(2), 1690–1698.
1714 <https://doi.org/10.1016/j.neuroimage.2009.10.017>
- 1715 Wiecki, T. V., Sofer, I., & Frank, M. J. (2013). HDDM: Hierarchical bayesian estimation of
1716 the drift-diffusion model in python. *Frontiers in Neuroinformatics*, 7.
1717 <https://doi.org/10.3389/fninf.2013.00014>
- 1718 Willis, J., & Todorov, A. (2006). First impressions: Making up your mind after a 100-ms
1719 exposure to a face. *Psychological Science*, 17(7), 592–598.
1720 <https://doi.org/10.1111/j.1467-9280.2006.01750.x>
- 1721 Xiao, Y. J., Coppin, G., & Bavel, J. J. V. (2016). Perceiving the world through
1722 group-colored glasses: A perceptual model of intergroup relations. *Psychological*
1723 *Inquiry*, 27(4), 255–274. <https://doi.org/10.1080/1047840X.2016.1199221>

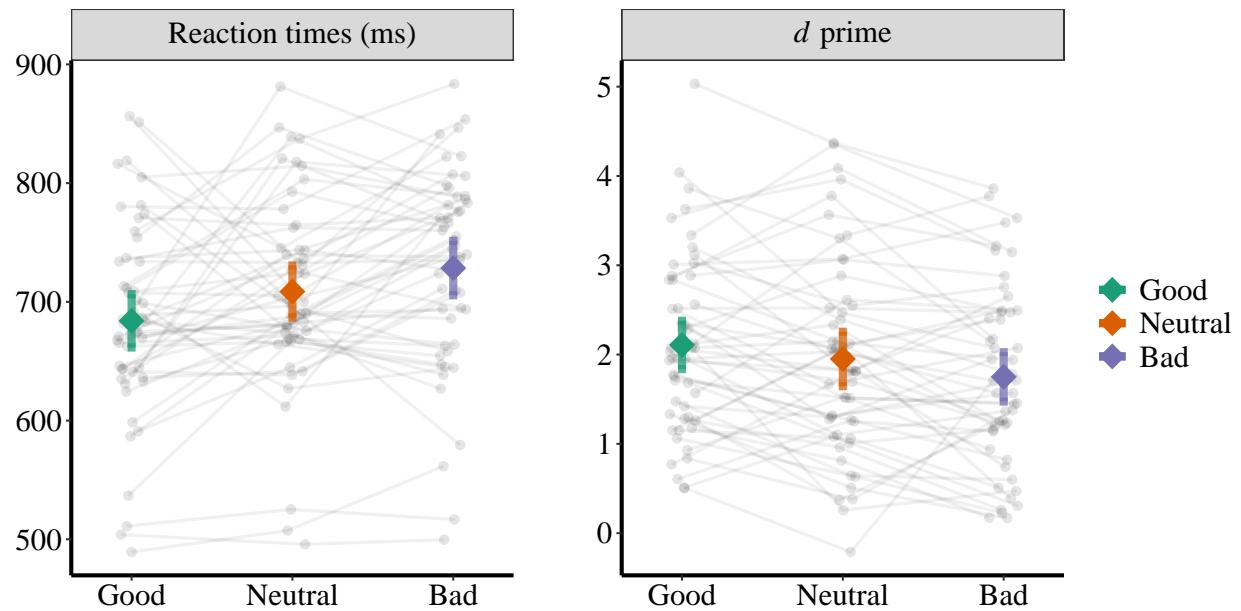


Figure 1. RT and d prime of Experiment 1a.

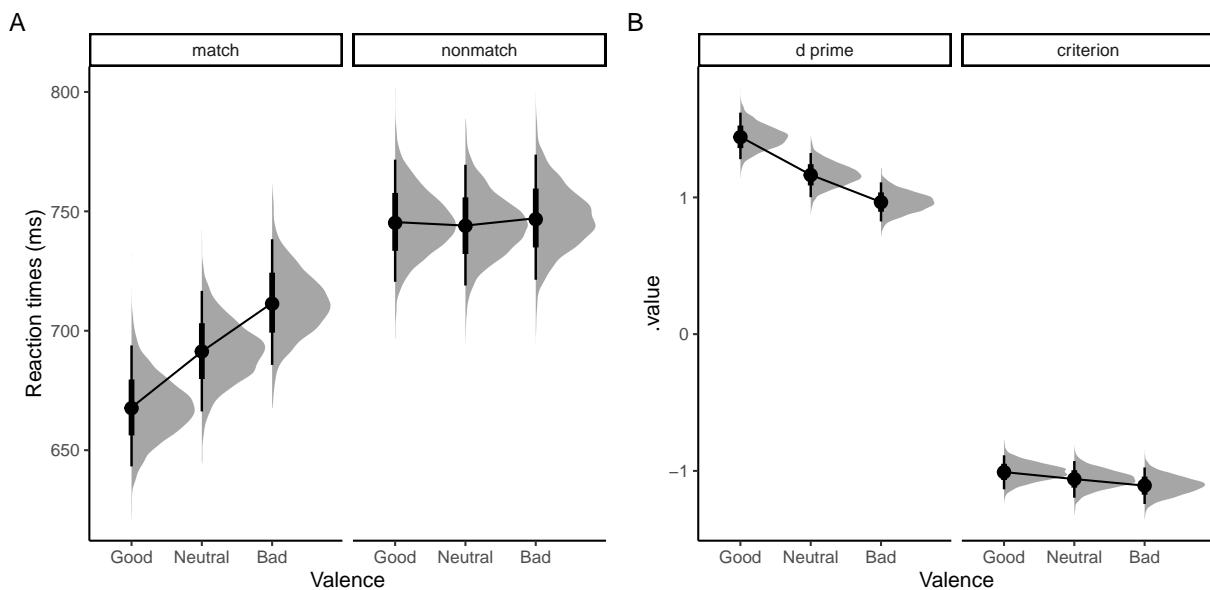


Figure 2. Exp1a: Results of Bayesian GLM analysis.

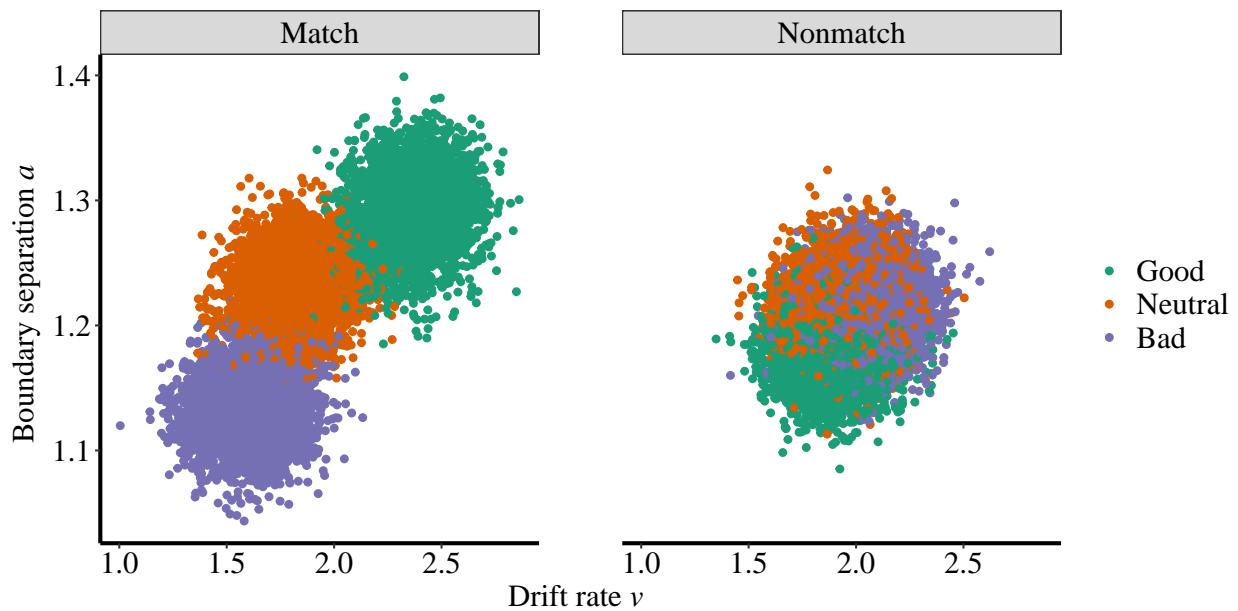


Figure 3. Exp1a: Results of HDDM.

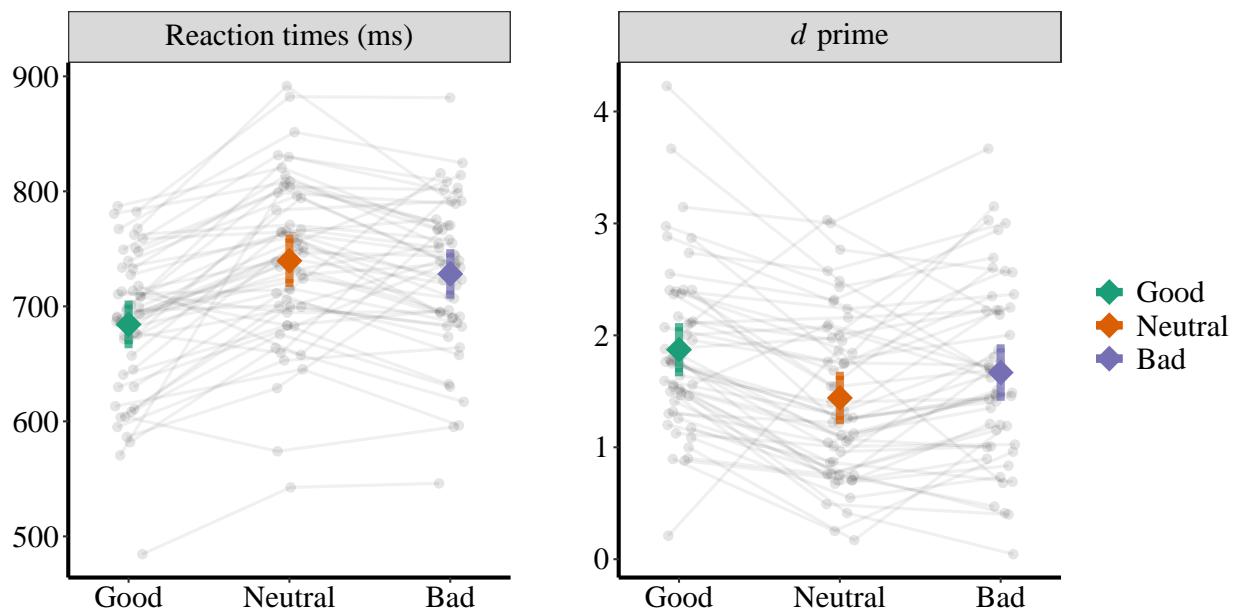


Figure 4. RT and d' of Experiment 1b.

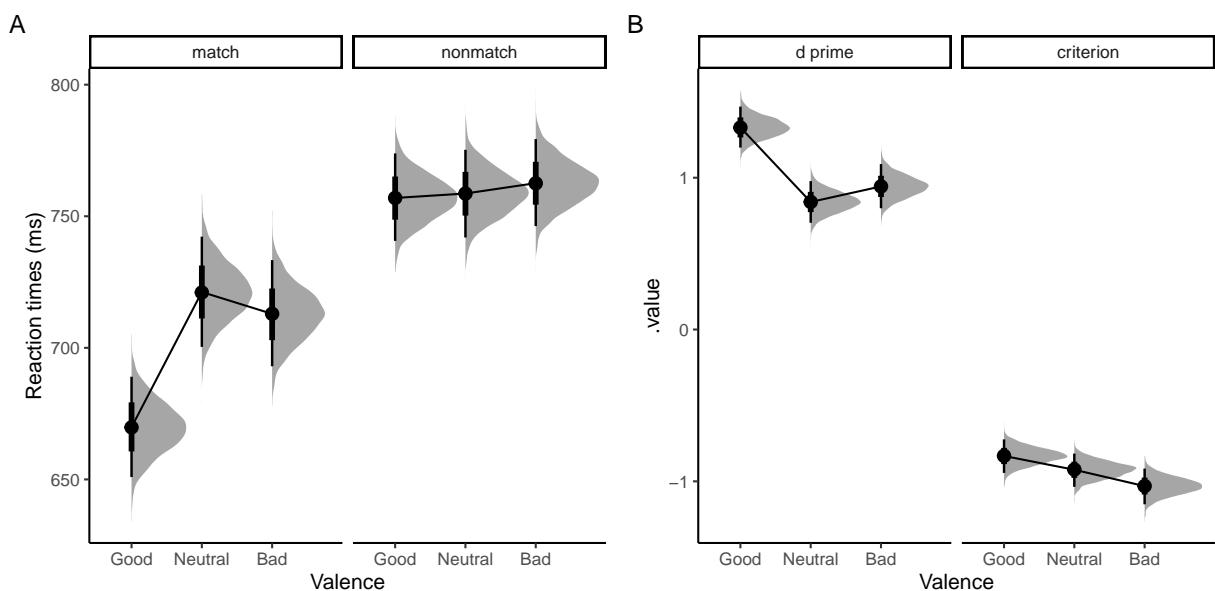


Figure 5. Exp1b: Results of Bayesian GLM analysis.

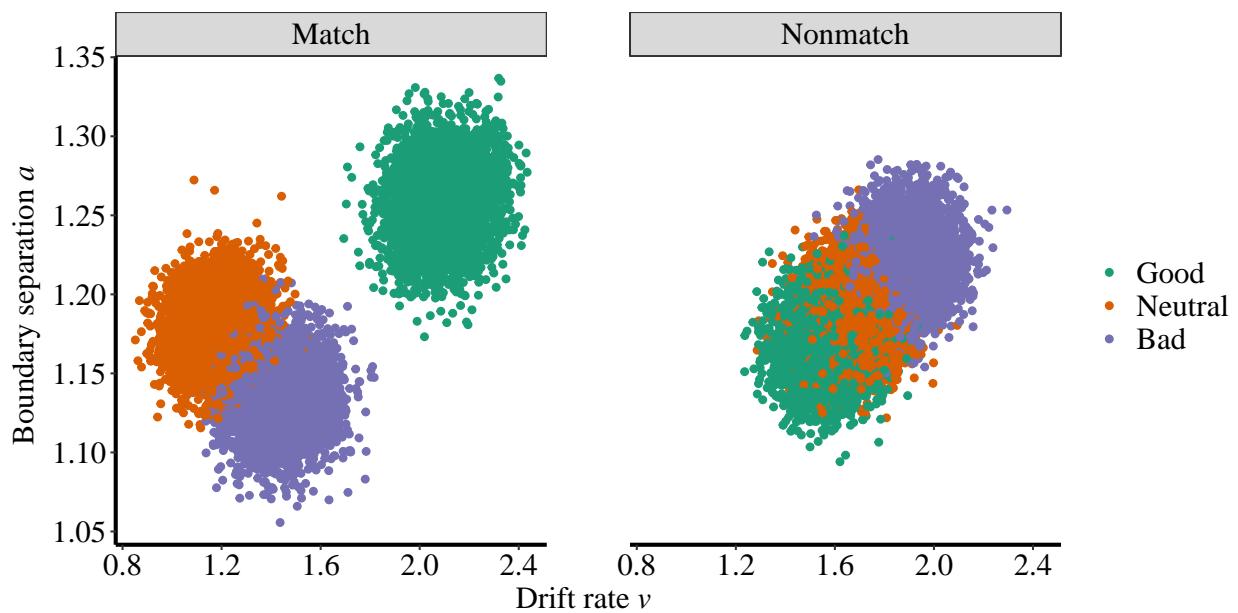


Figure 6. Exp1b: Results of HDDM.

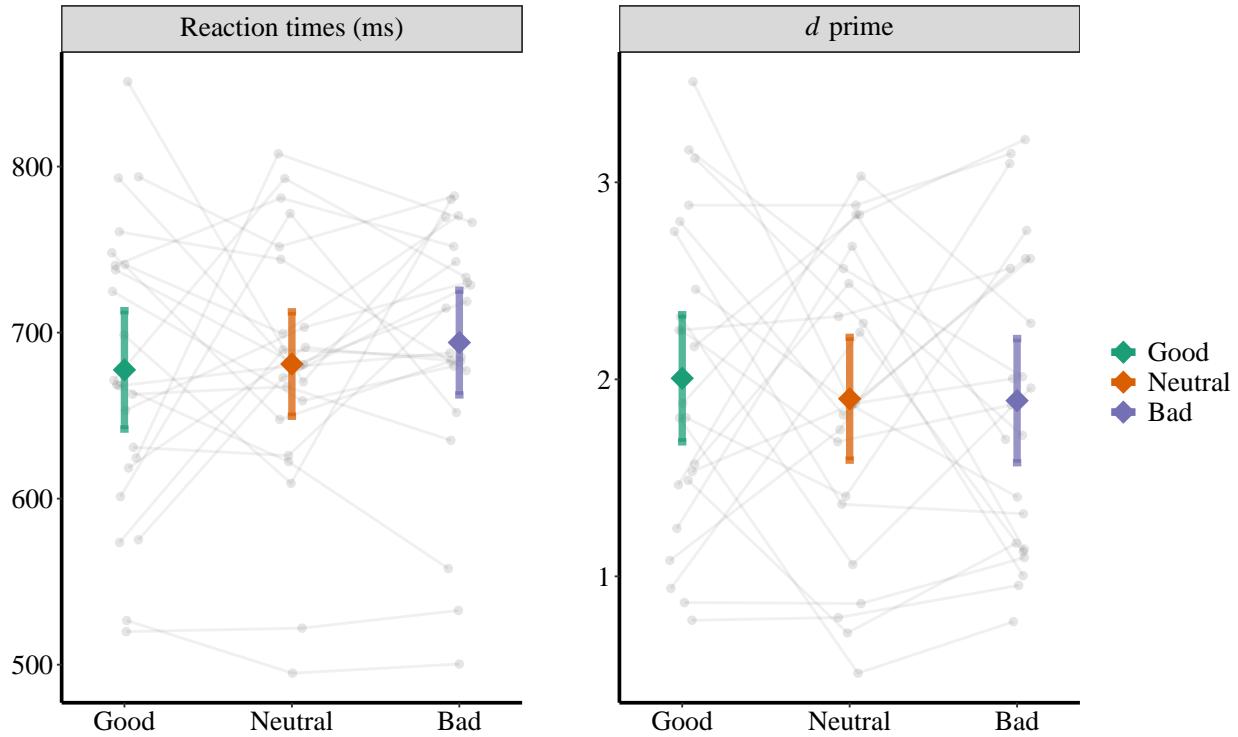


Figure 7. RT and d' prime of Experiment 1c.

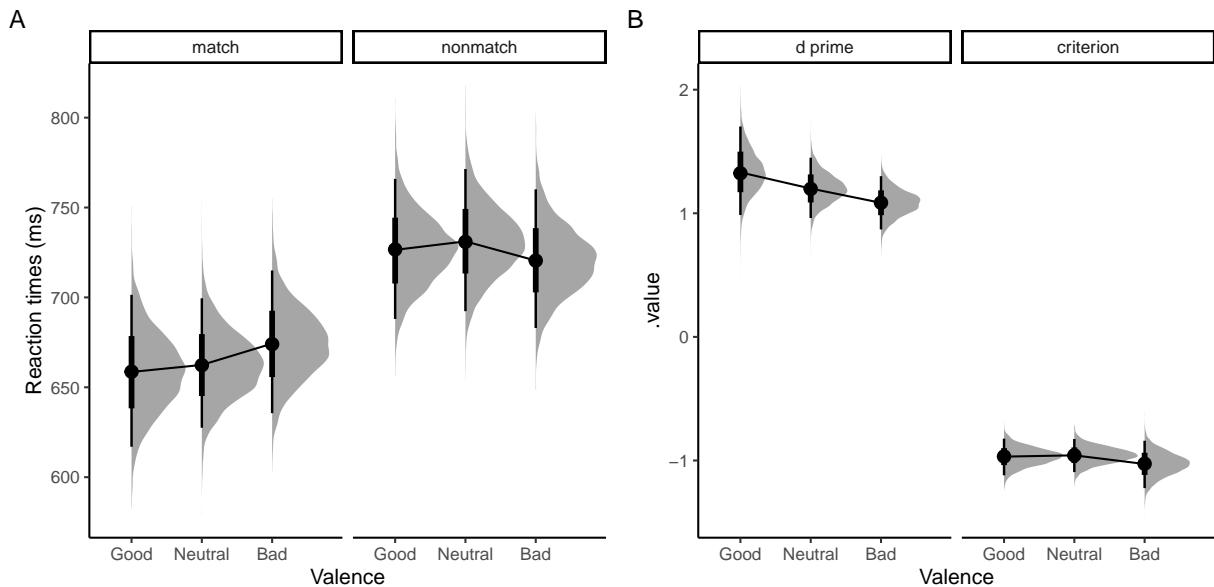


Figure 8. Exp1c: Results of Bayesian GLM analysis.

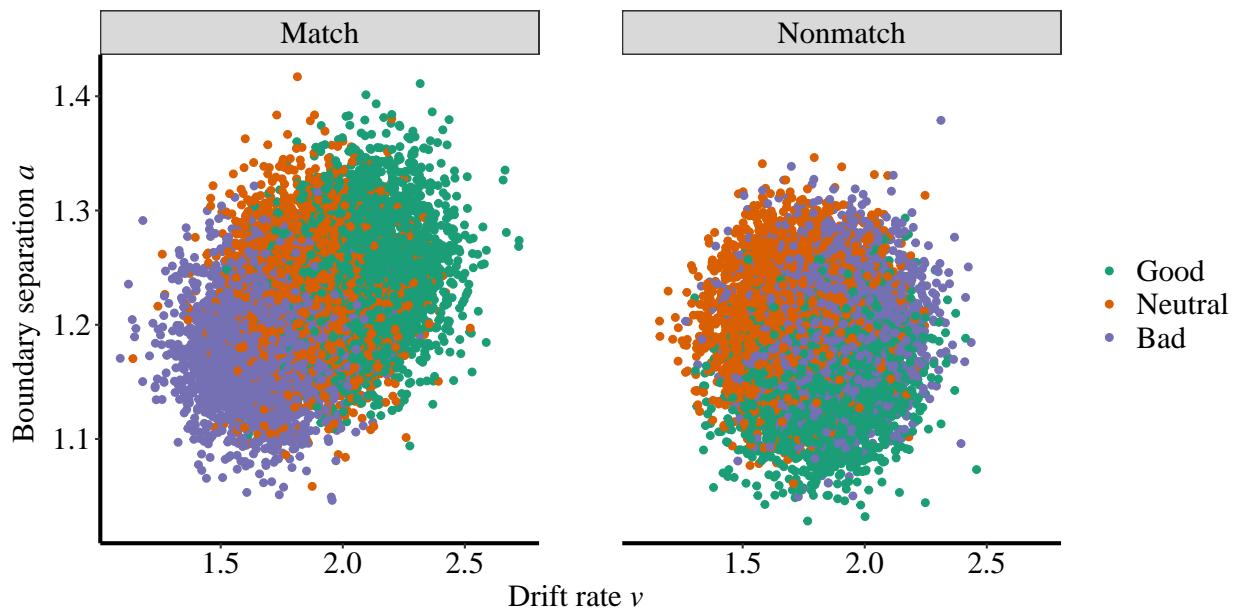


Figure 9. Exp1c: Results of HDDM.

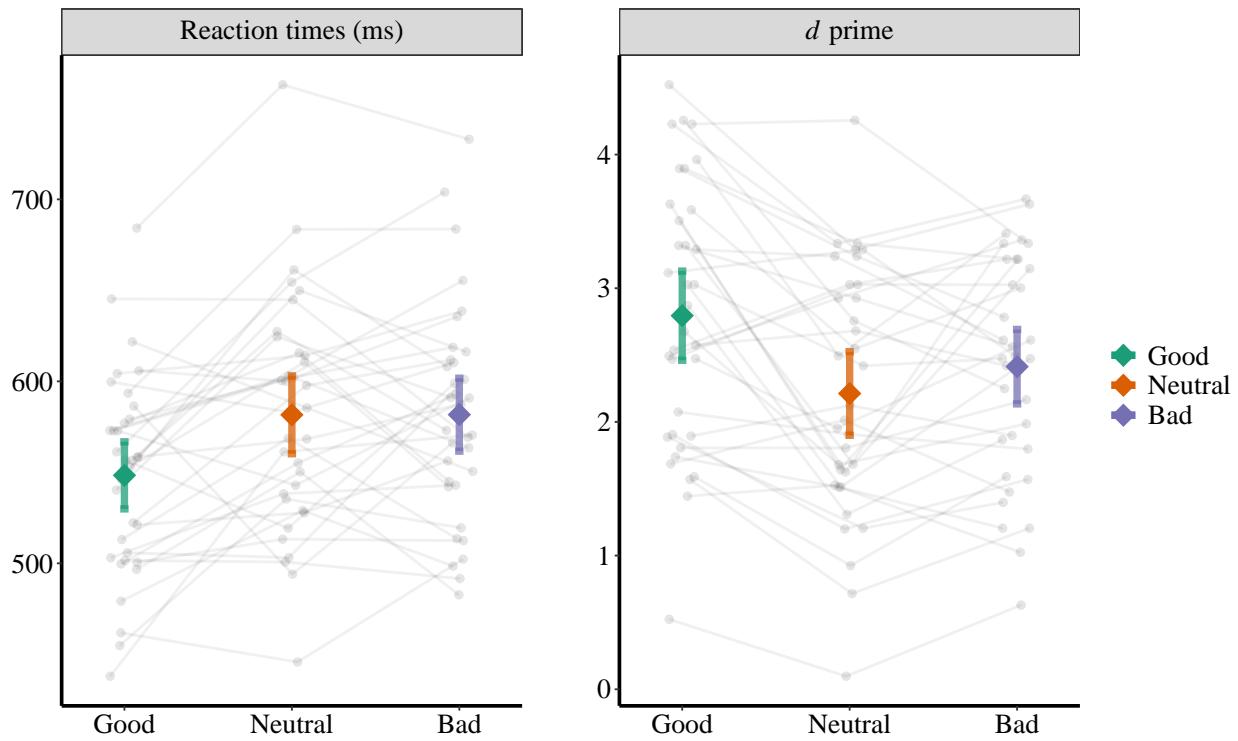


Figure 10. RT and d' of Experiment 2.

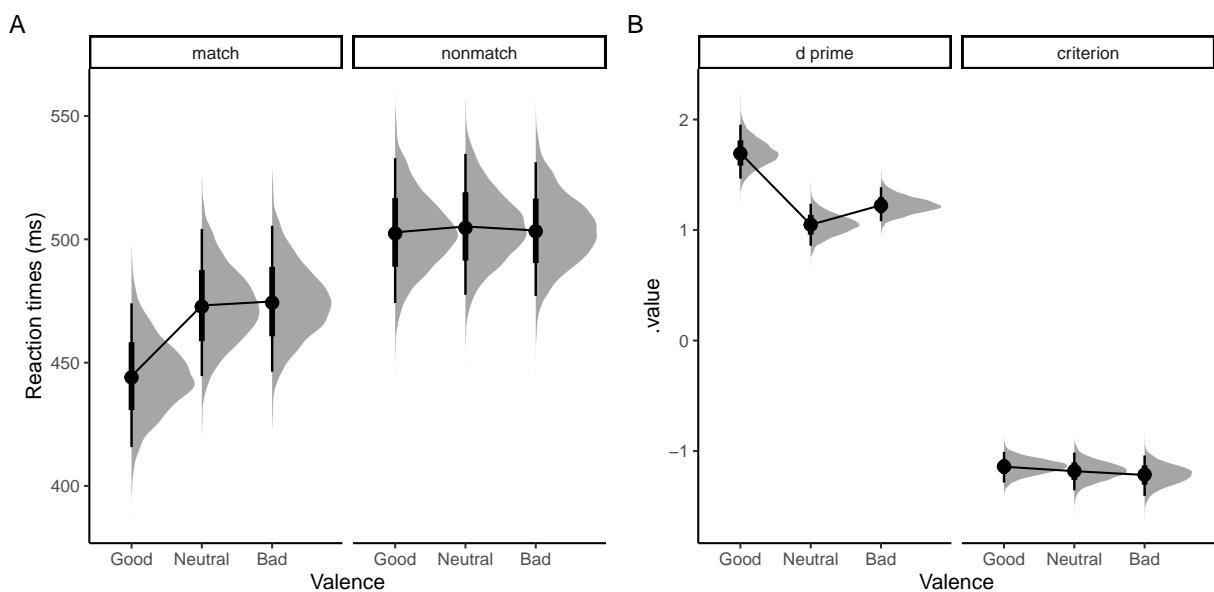


Figure 11. Exp2: Results of Bayesian GLM analysis.

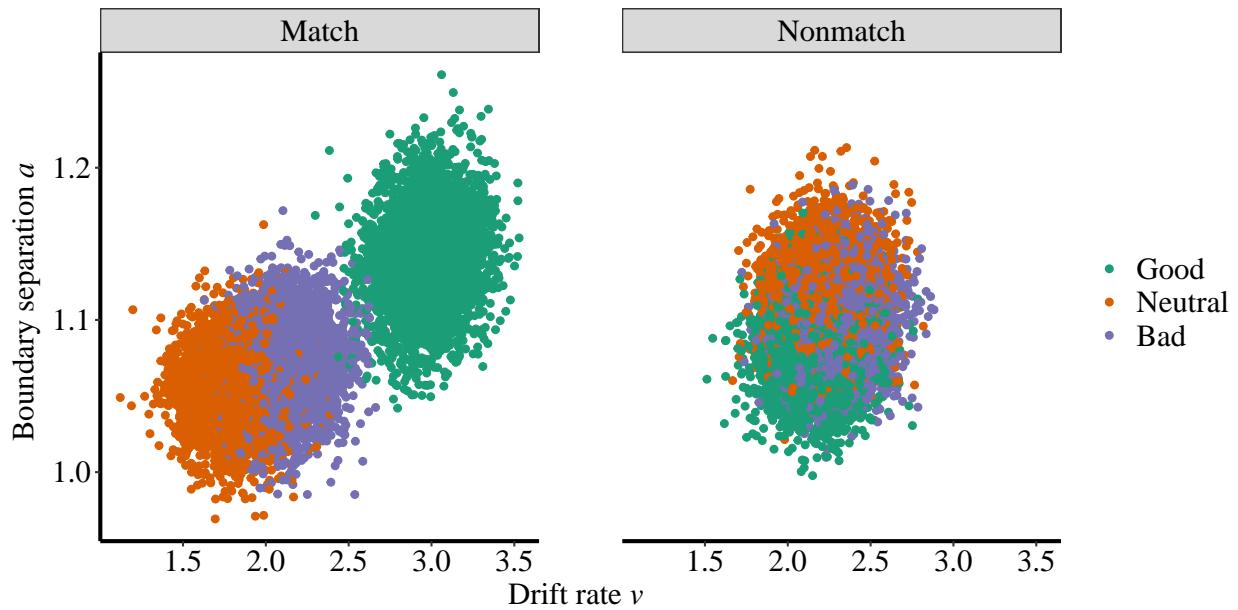


Figure 12. Exp2: Results of HDDM.

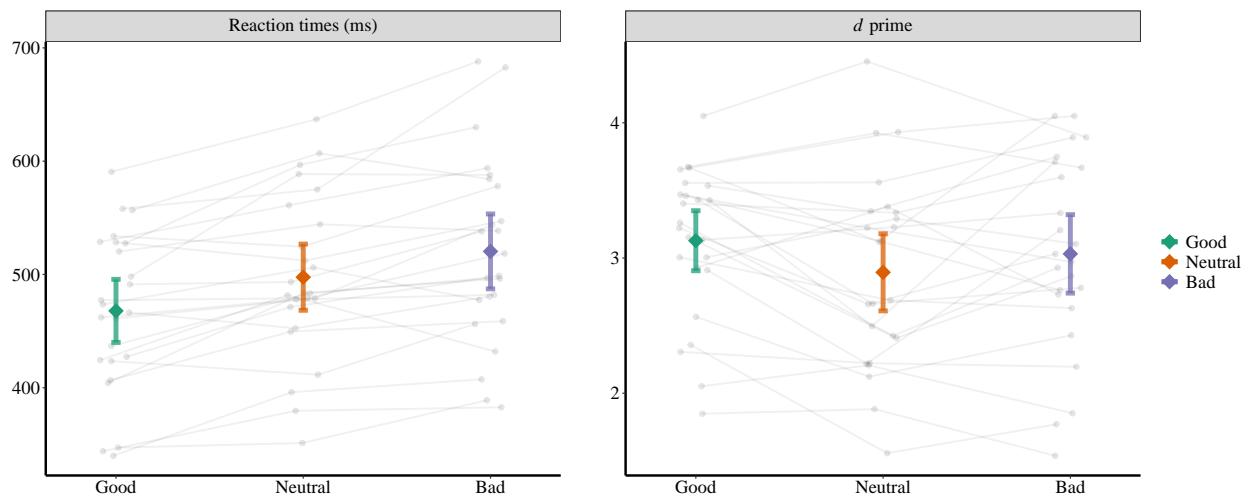


Figure 13. RT and d' prime of Experiment 6a.

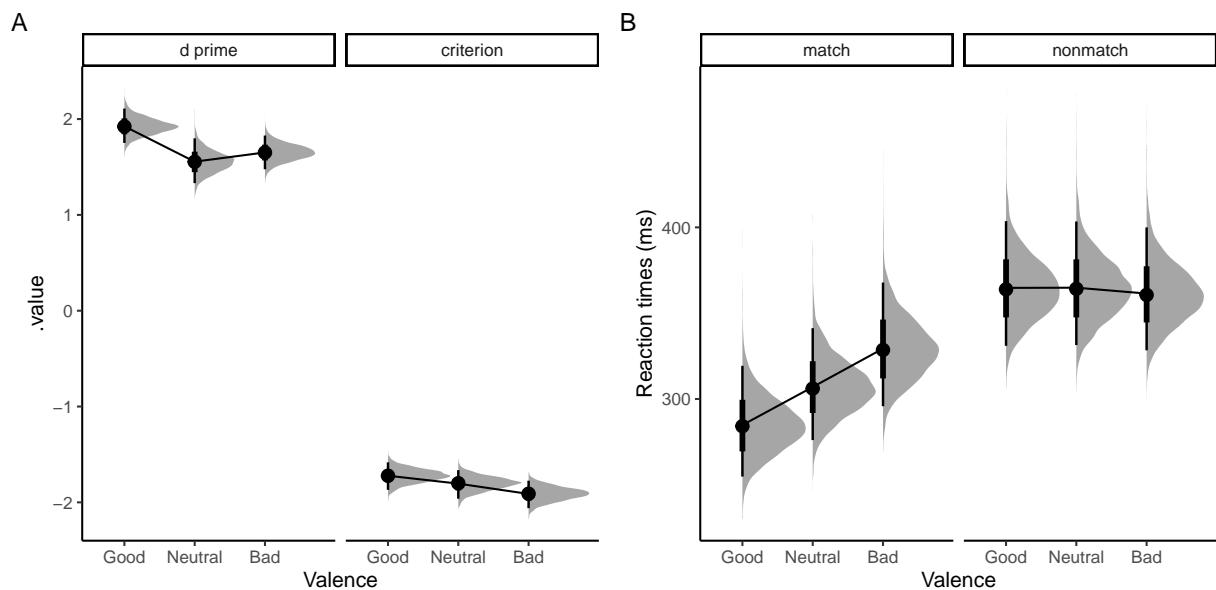


Figure 14. Exp6a: Results of Bayesian GLM analysis.

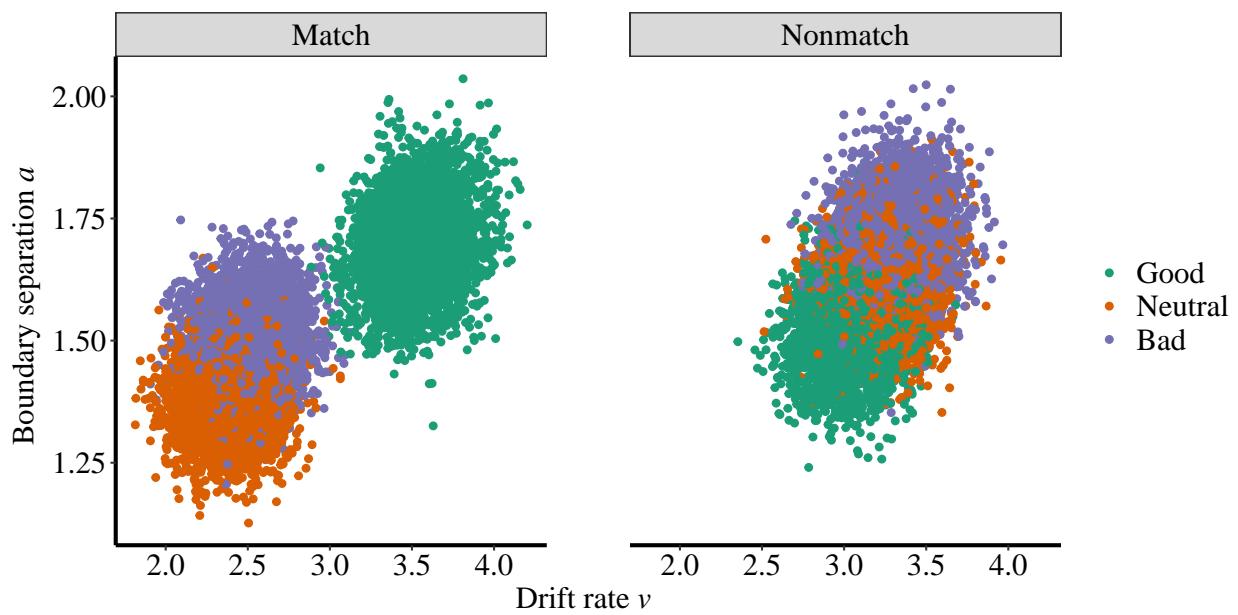


Figure 15. exp6a: Results of HDDM.

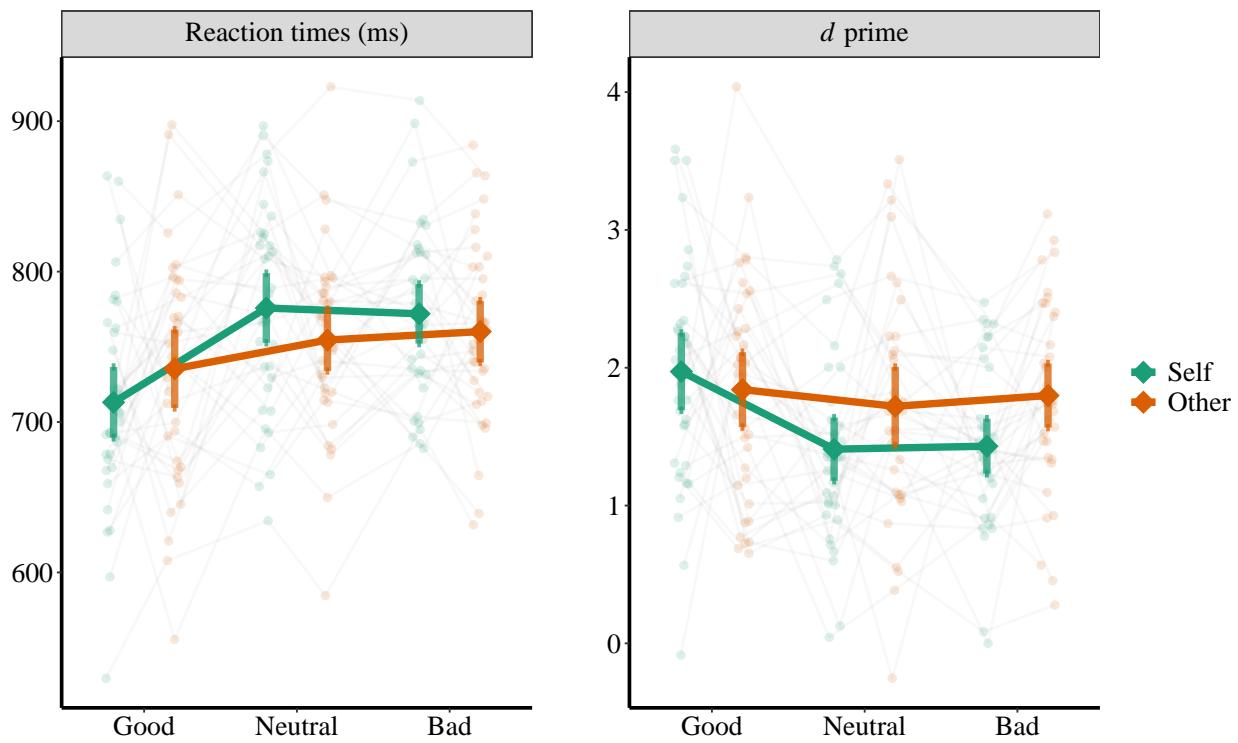


Figure 16. RT and d prime of Experiment 3a.

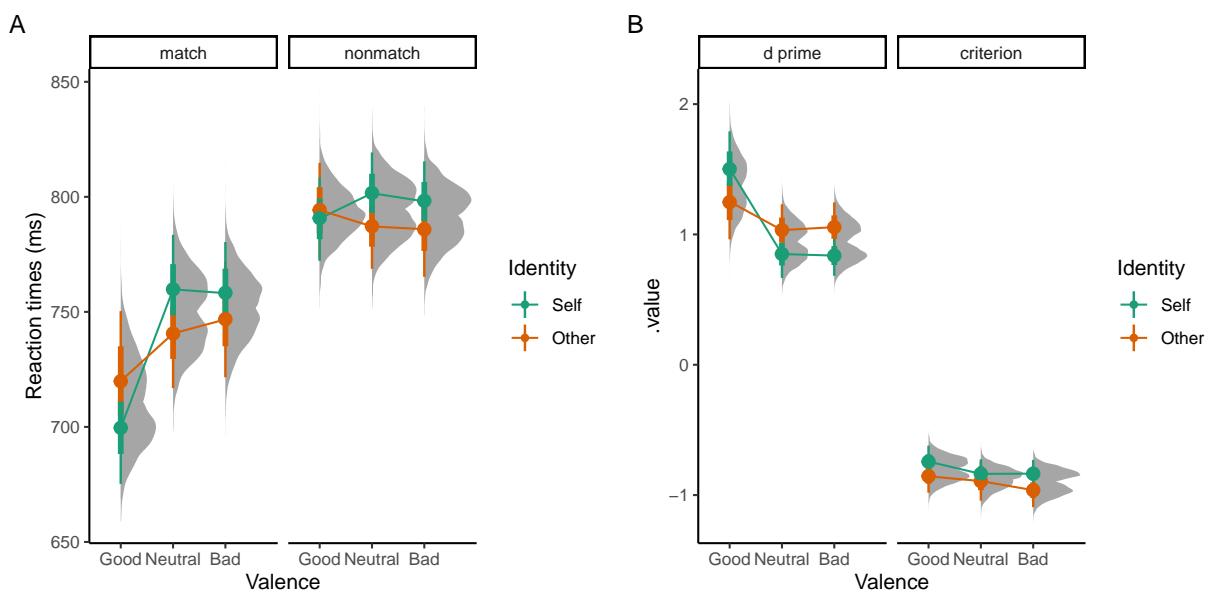


Figure 17. Exp3a: Results of Bayesian GLM analysis.

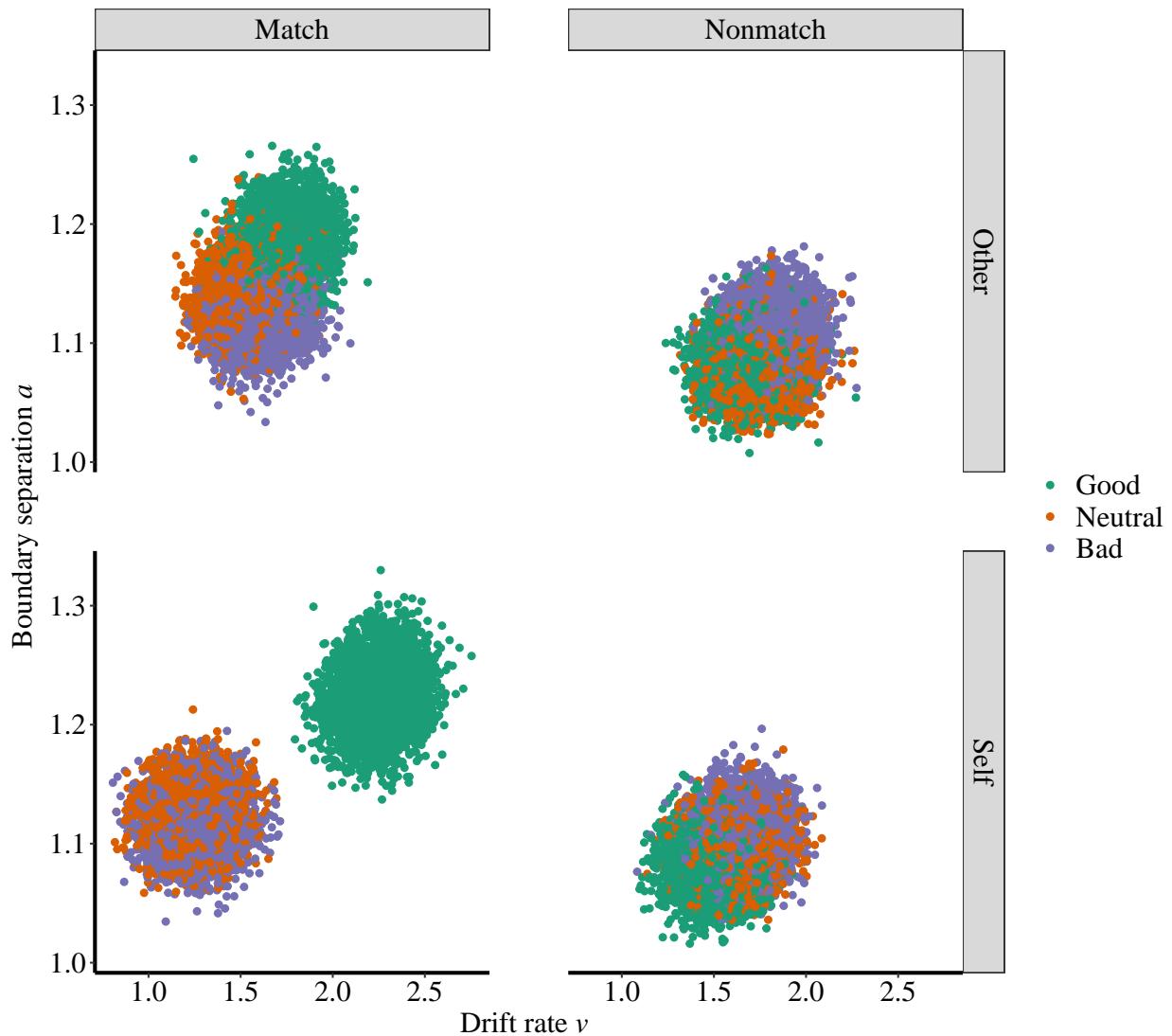


Figure 18. Exp3a: Results of HDDM.

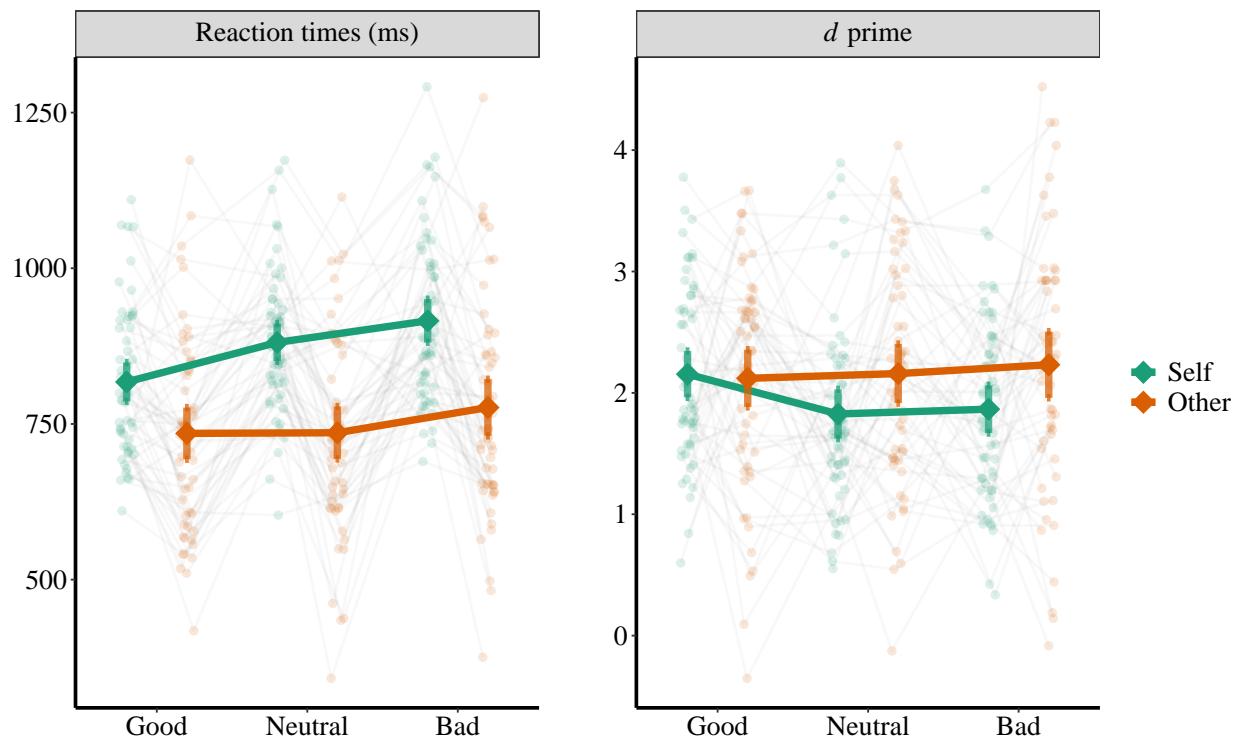


Figure 19. RT and d' prime of Experiment 3b.

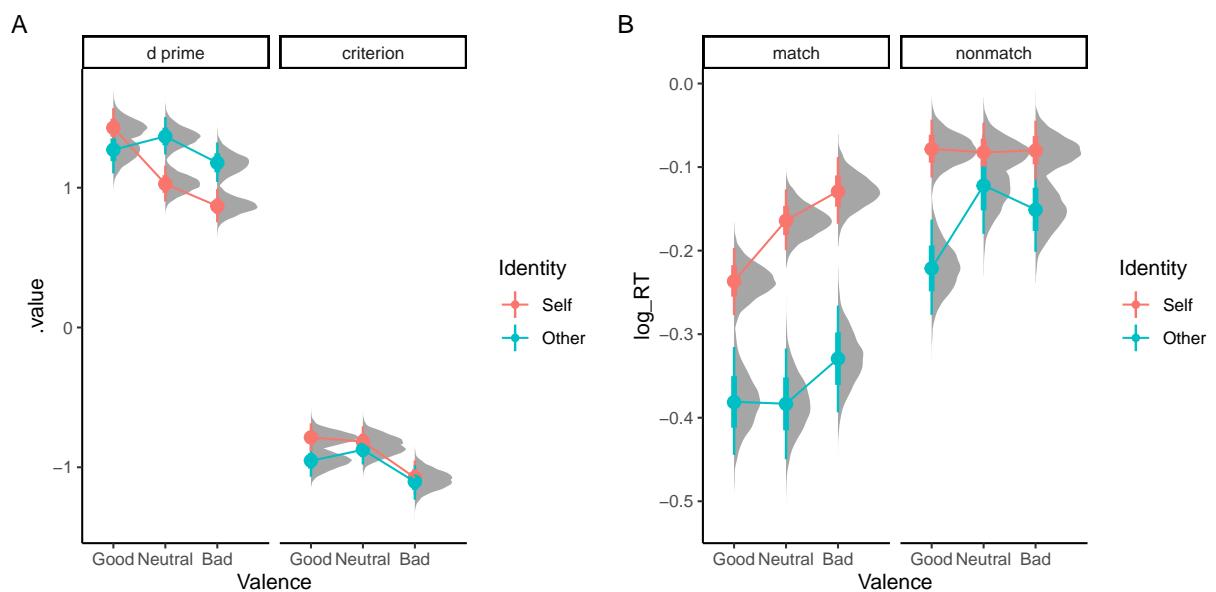


Figure 20. exp3b: Results of Bayesian GLM analysis.

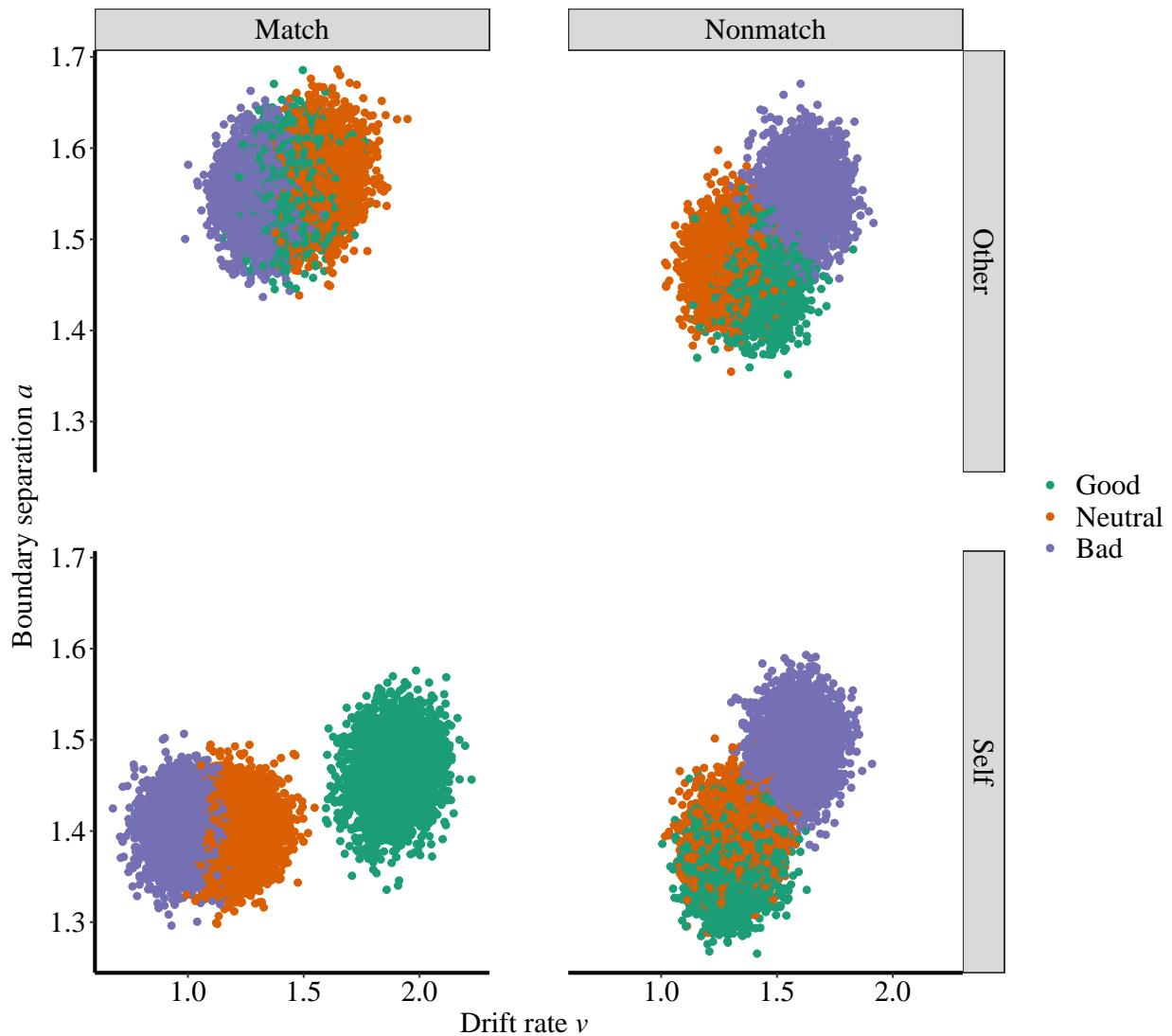


Figure 21. exp3b: Results of HDDM.

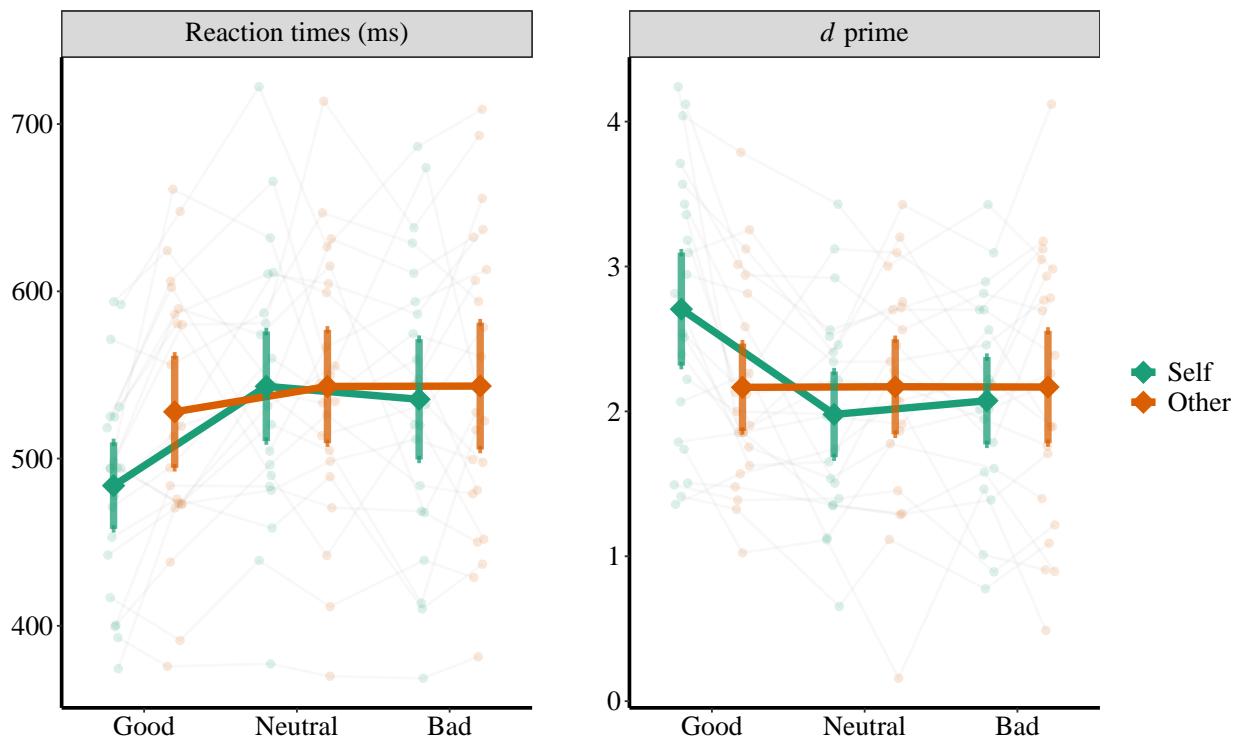


Figure 22. RT and d prime of Experiment 6b.

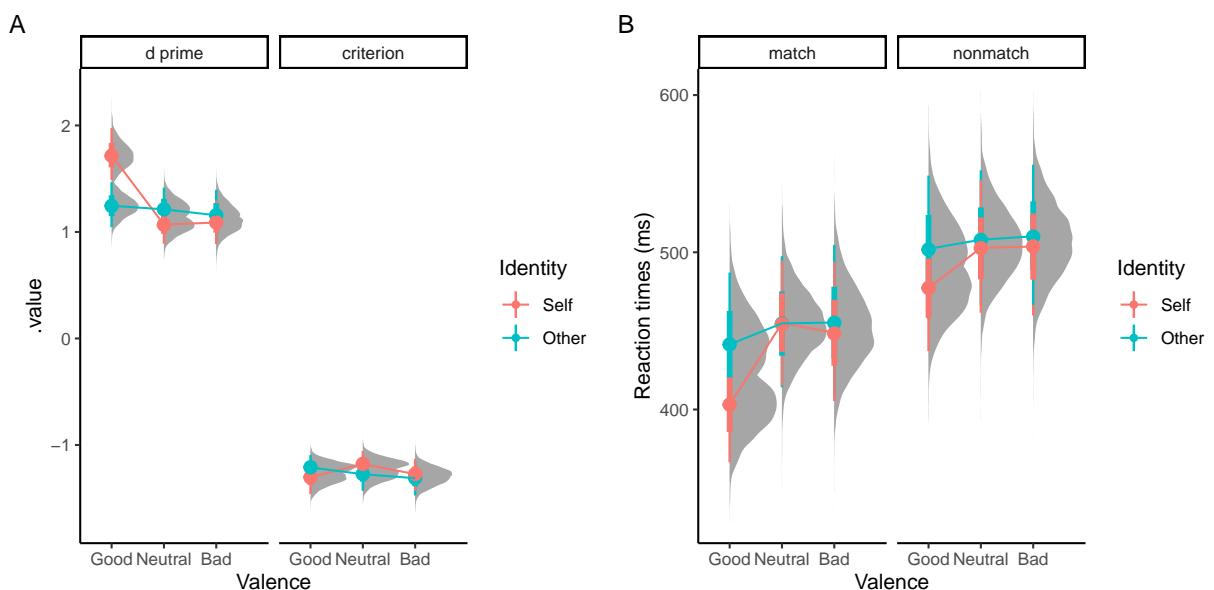


Figure 23. exp6b_d1: Results of Bayesian GLM analysis.

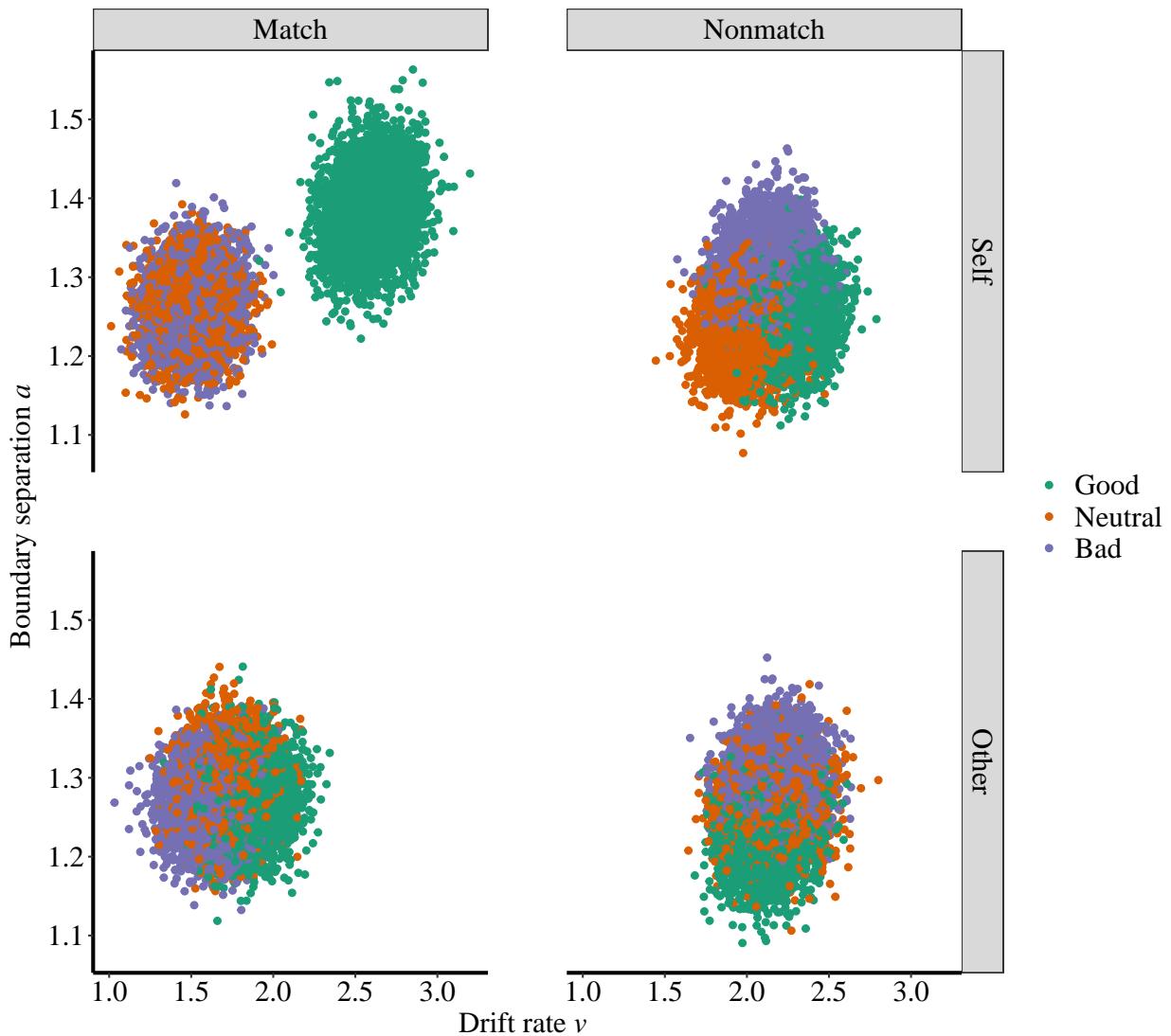


Figure 24. exp6b: Results of HDDM (Day 1).

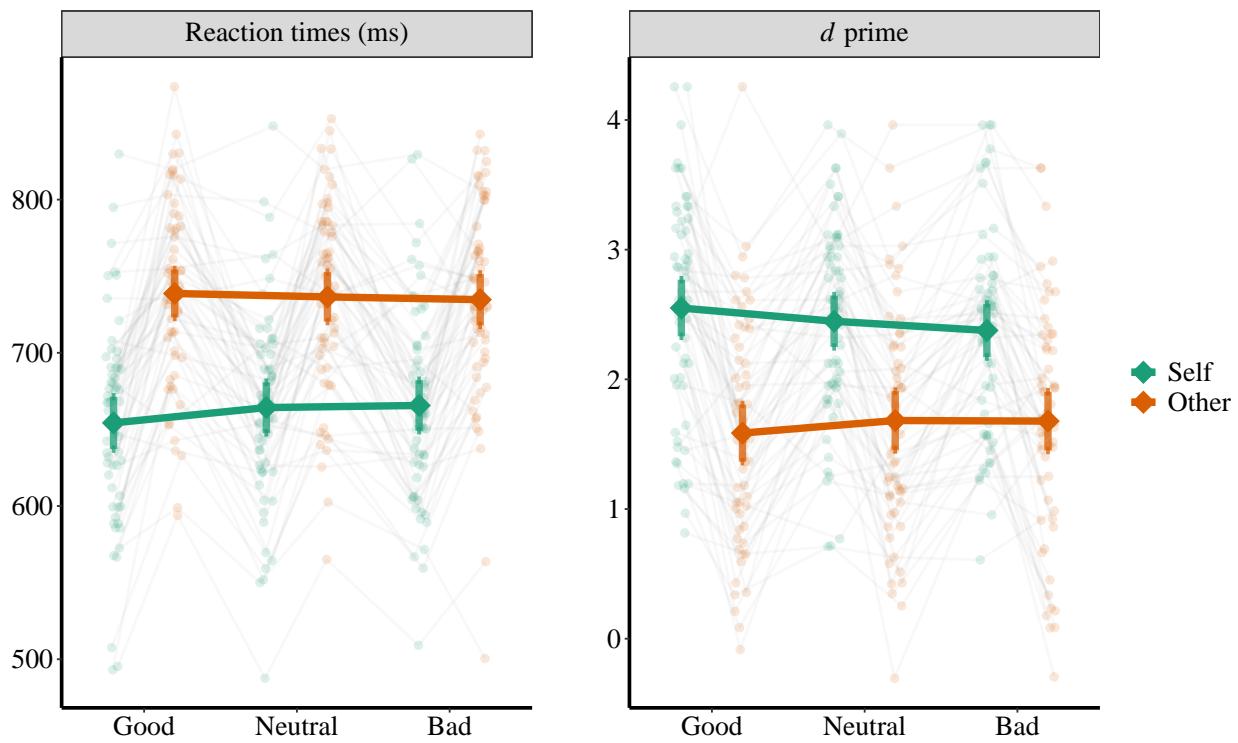
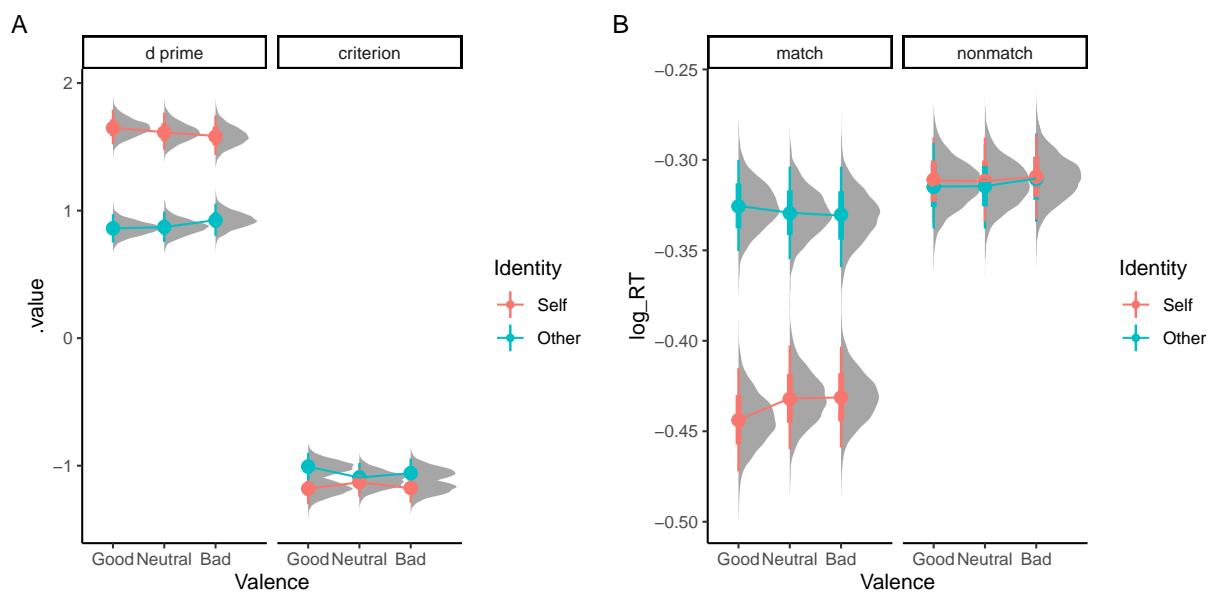
Figure 25. RT and d' of Experiment 4a.

Figure 26. exp4a: Results of Bayesian GLM analysis.

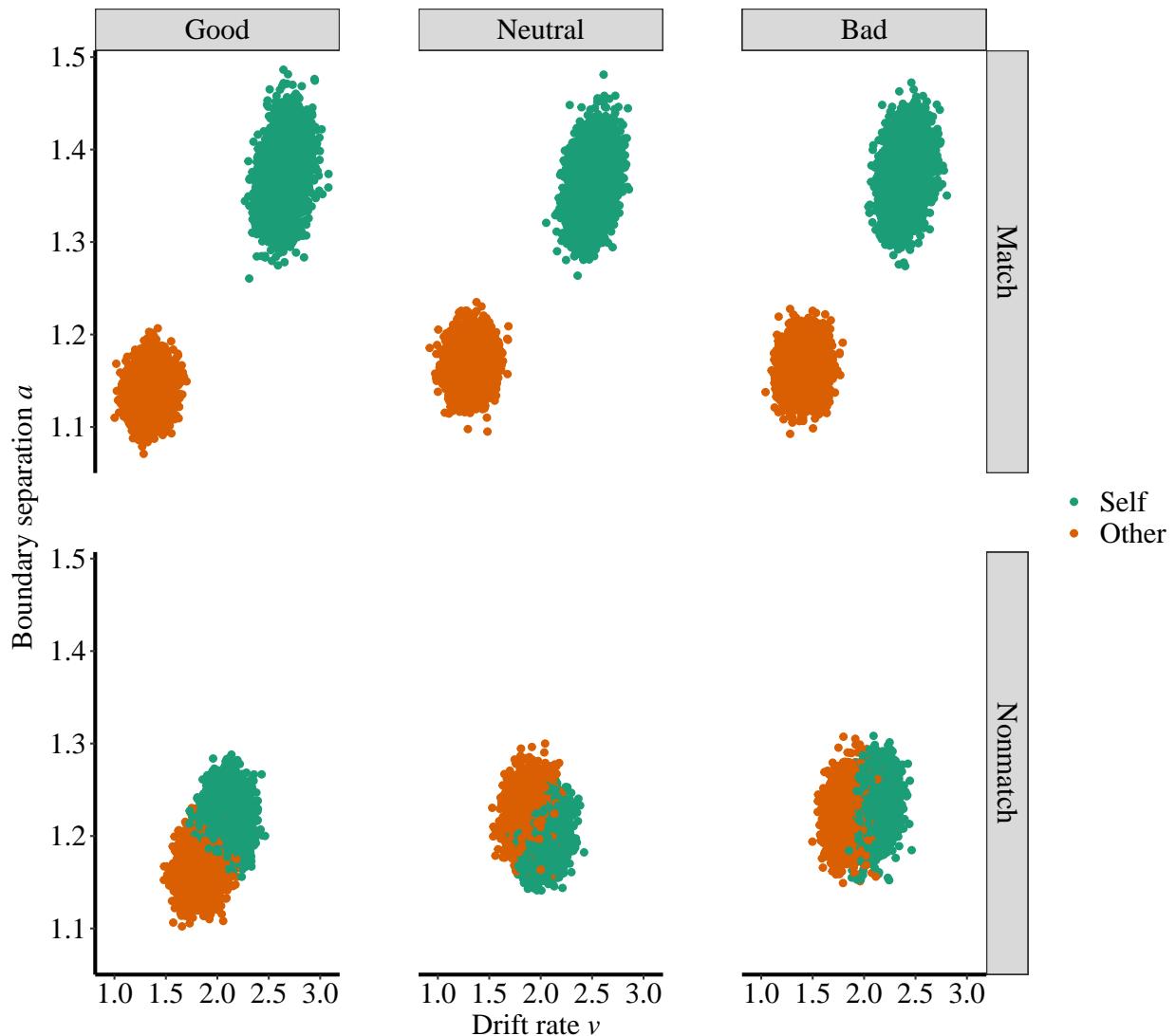


Figure 27. exp4a: Results of HDDM.

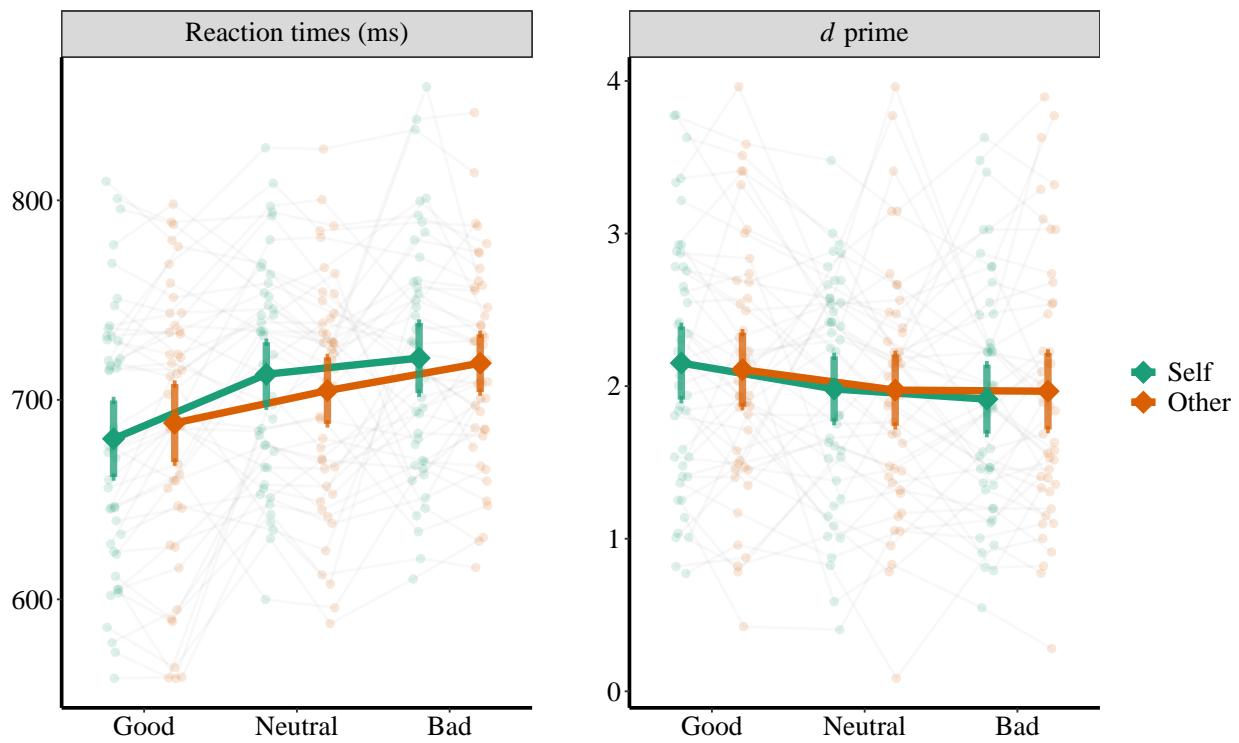


Figure 28. RT and d' of Experiment 4b.

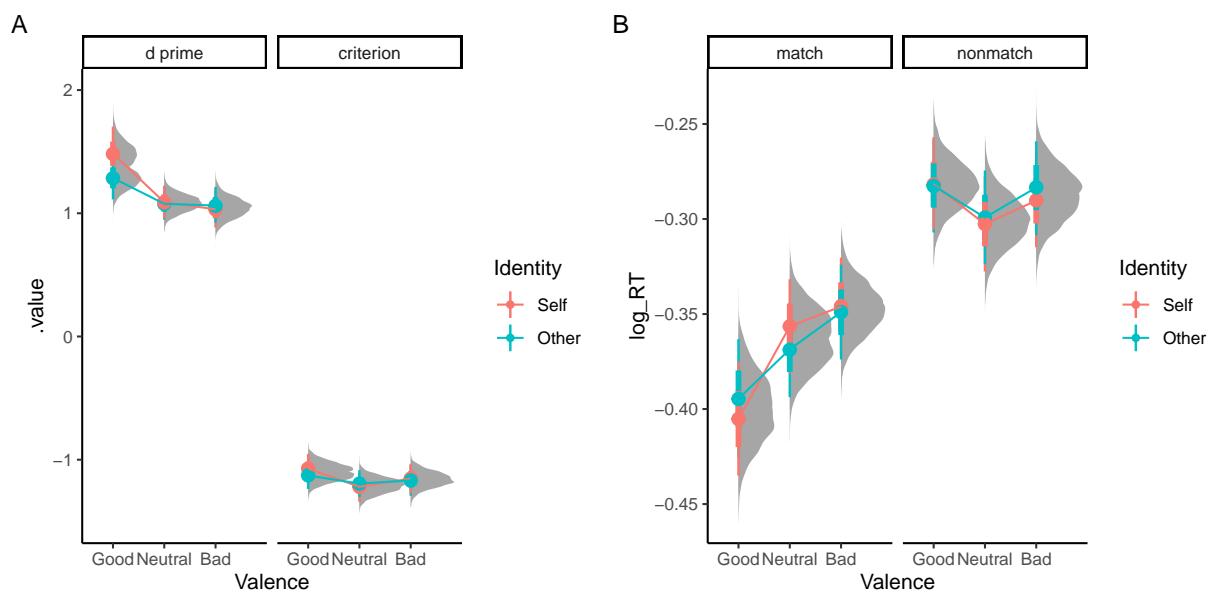


Figure 29. exp4b: Results of Bayesian GLM analysis.

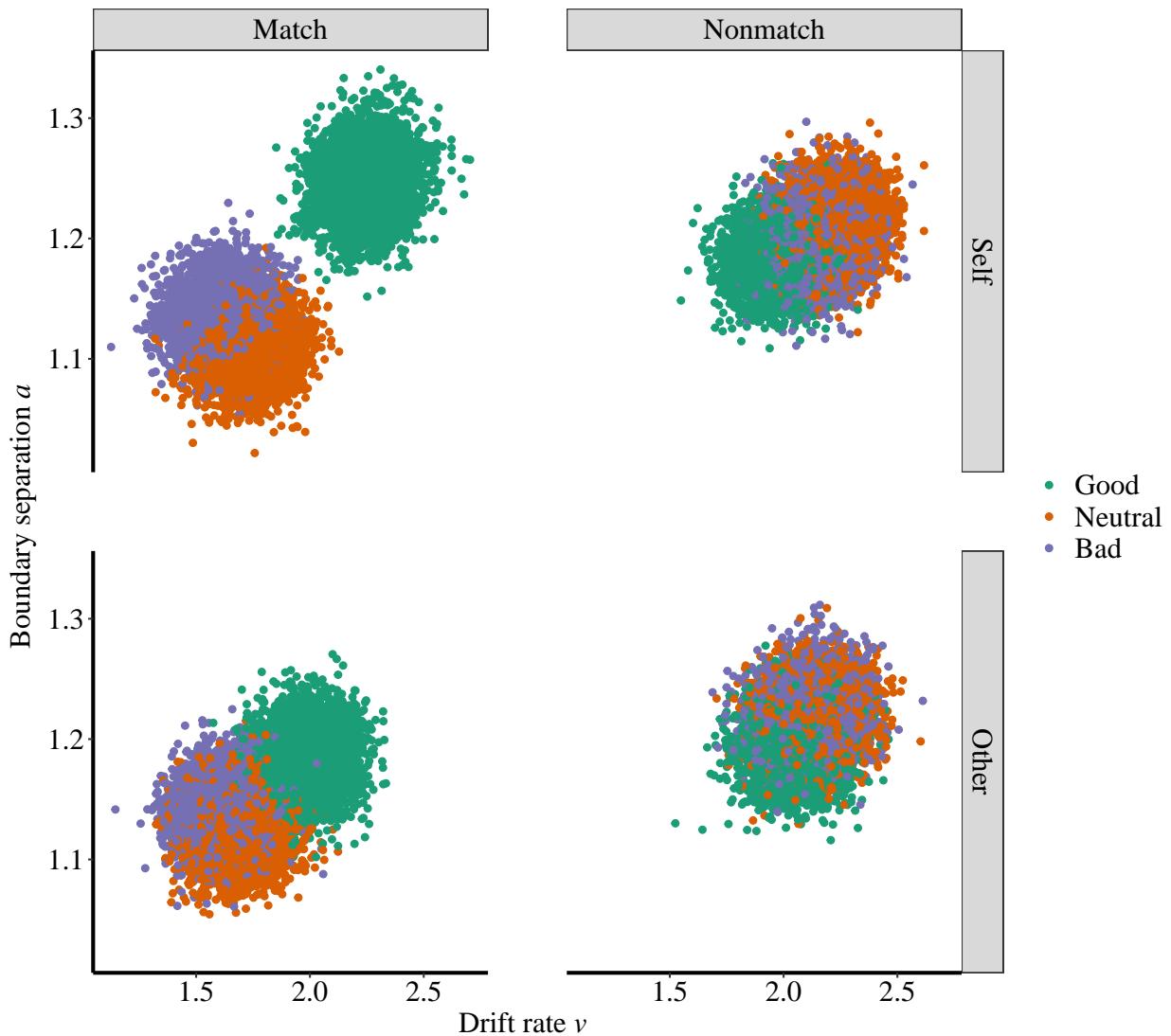


Figure 30. exp4b: Results of HDDM.

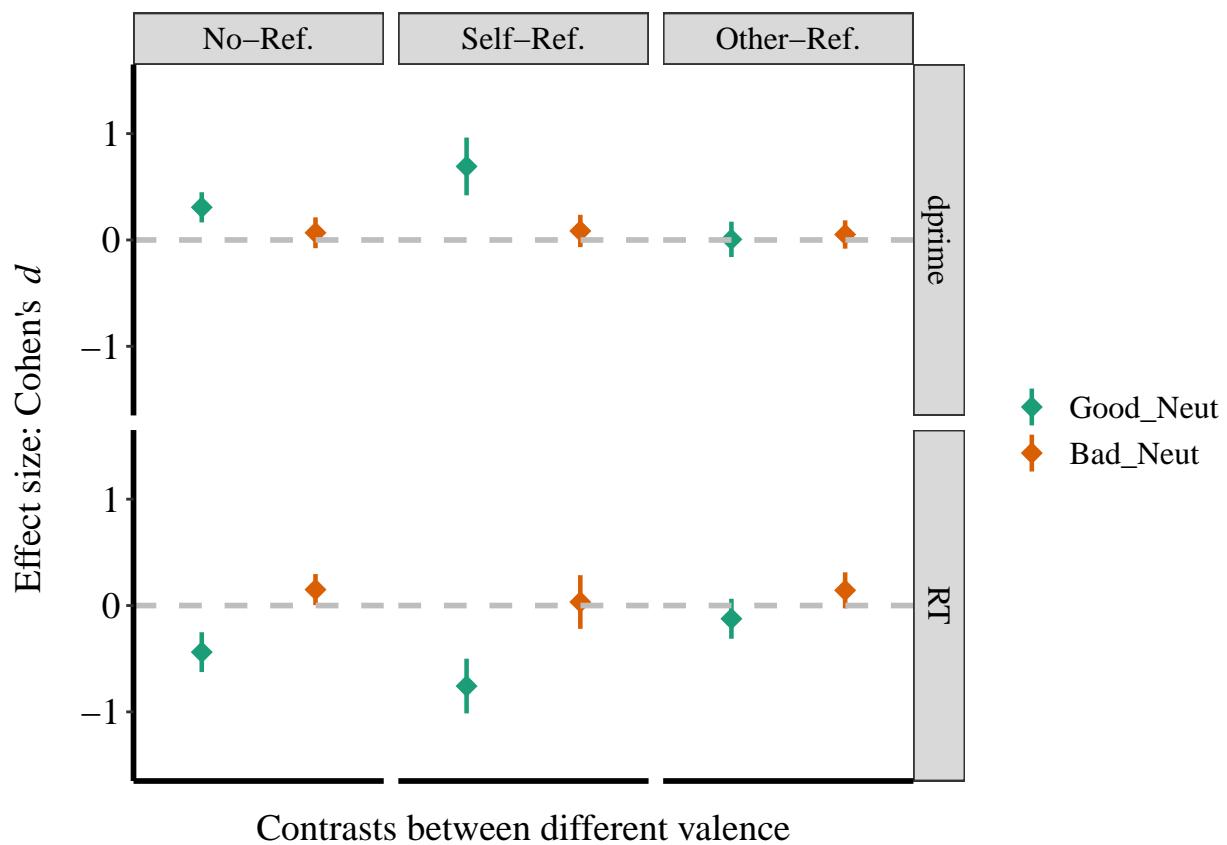


Figure 31. Effect size (Cohen's d) of Valence.

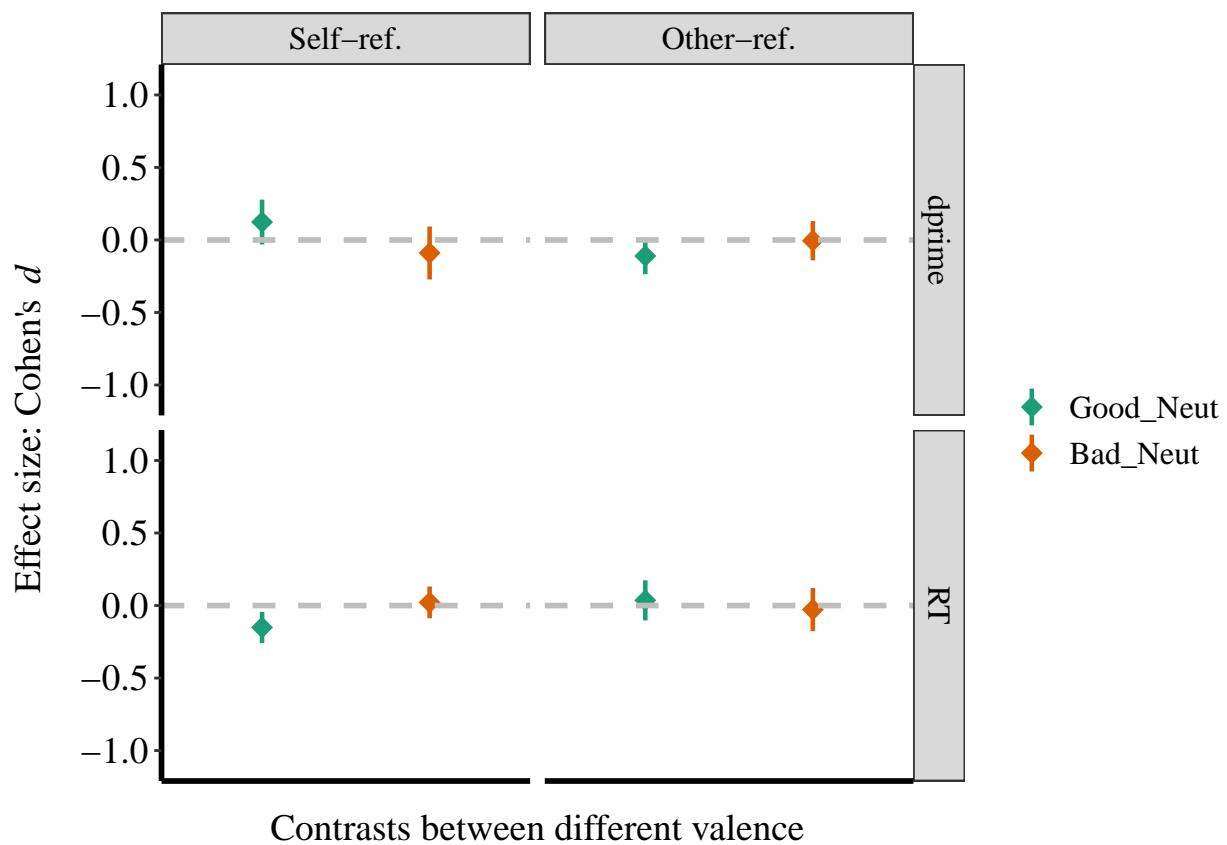


Figure 32. Effect size (Cohen's d) of Valence in Exp4a.

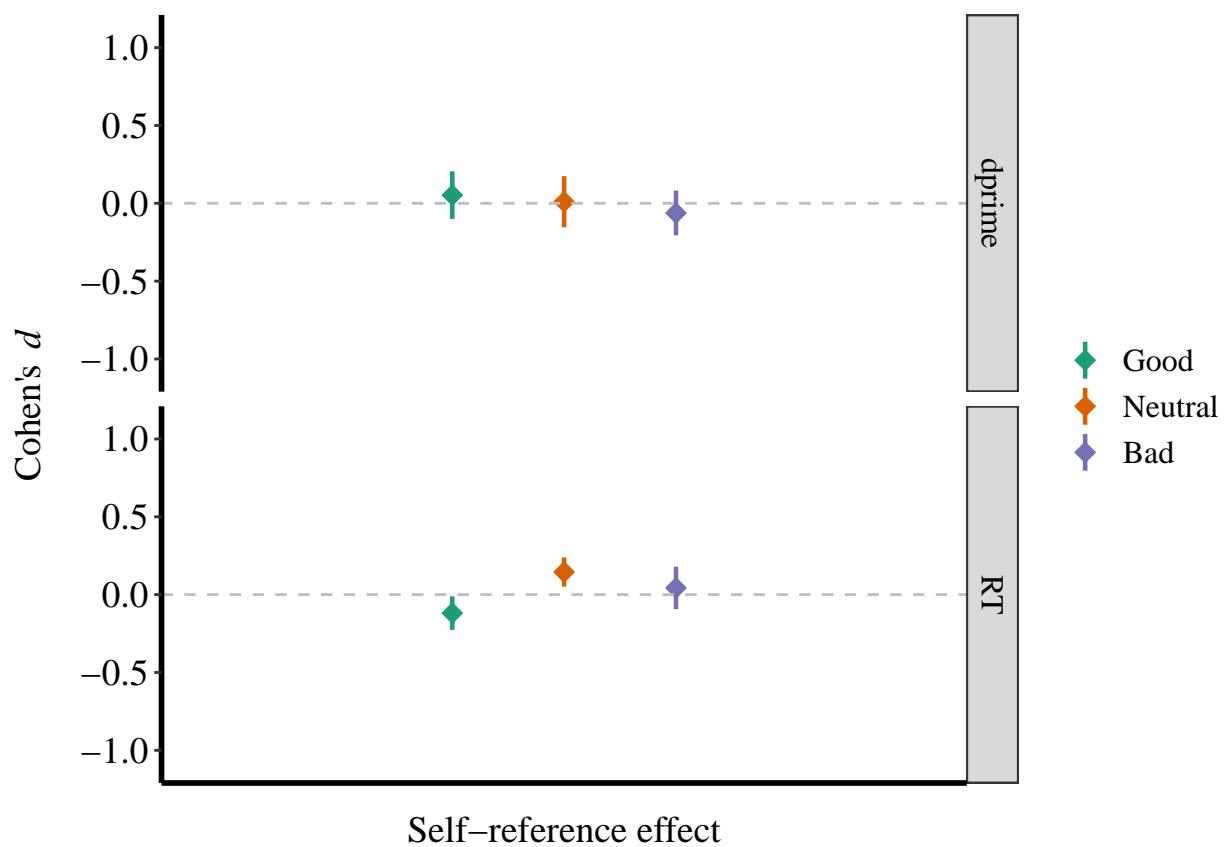


Figure 33. Effect size (Cohen's d) of Valence in Exp4b.

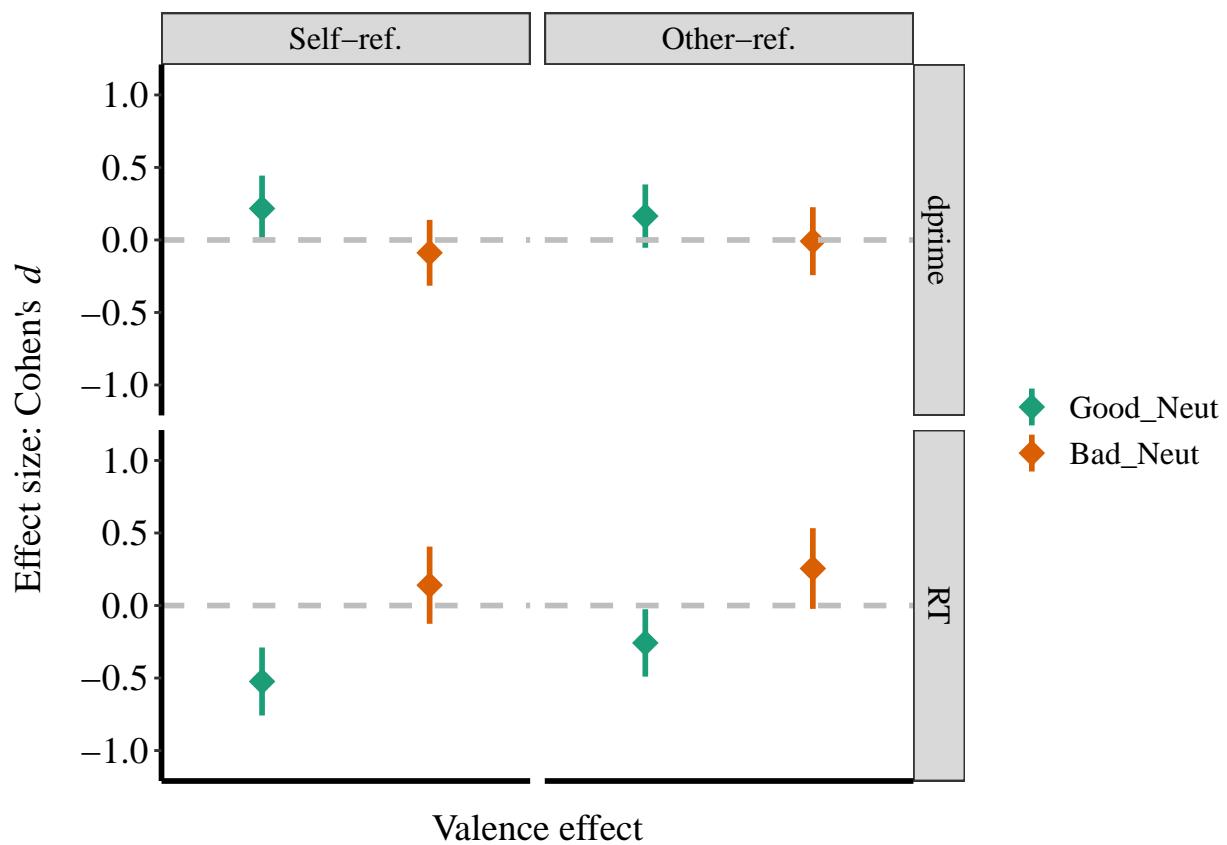


Figure 34. Effect size (Cohen's d) of Valence in Exp4b.

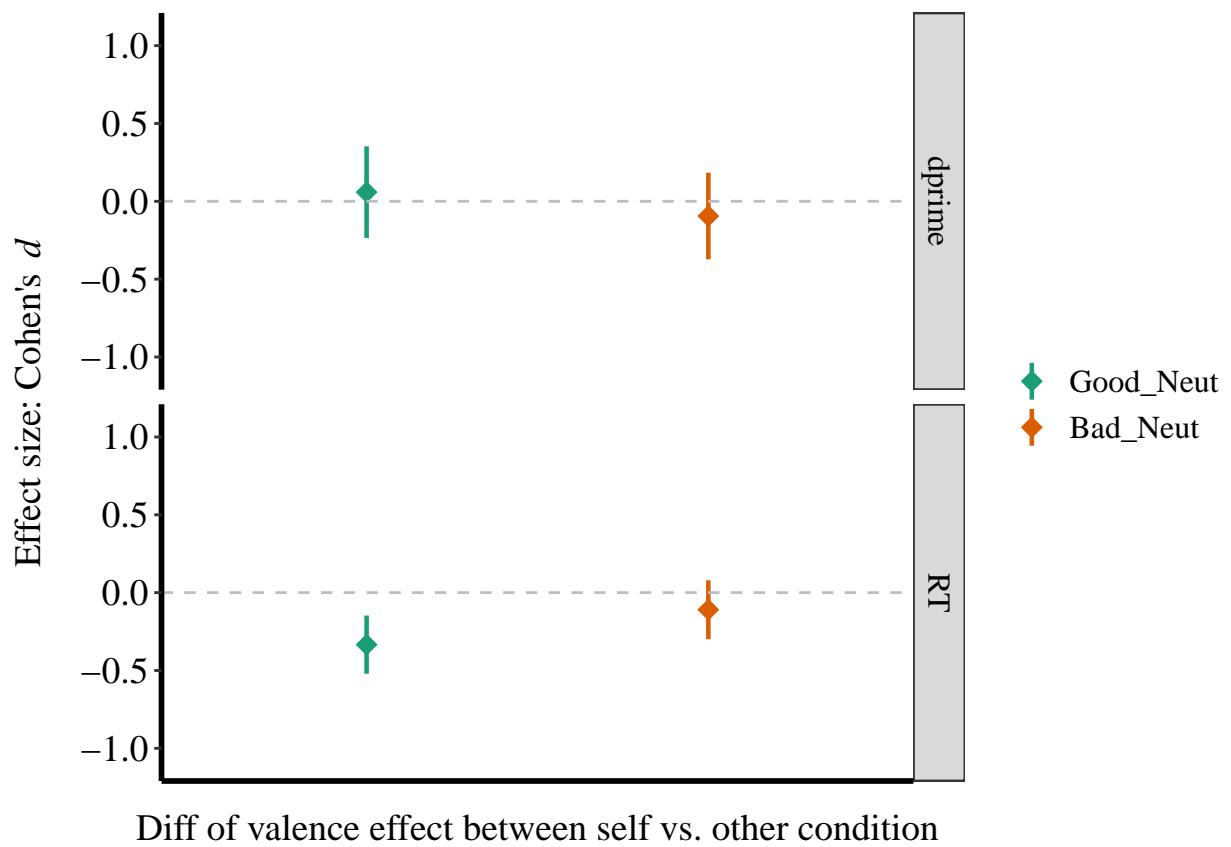


Figure 35. Effect size (Cohen's d) of Valence in Exp4b.

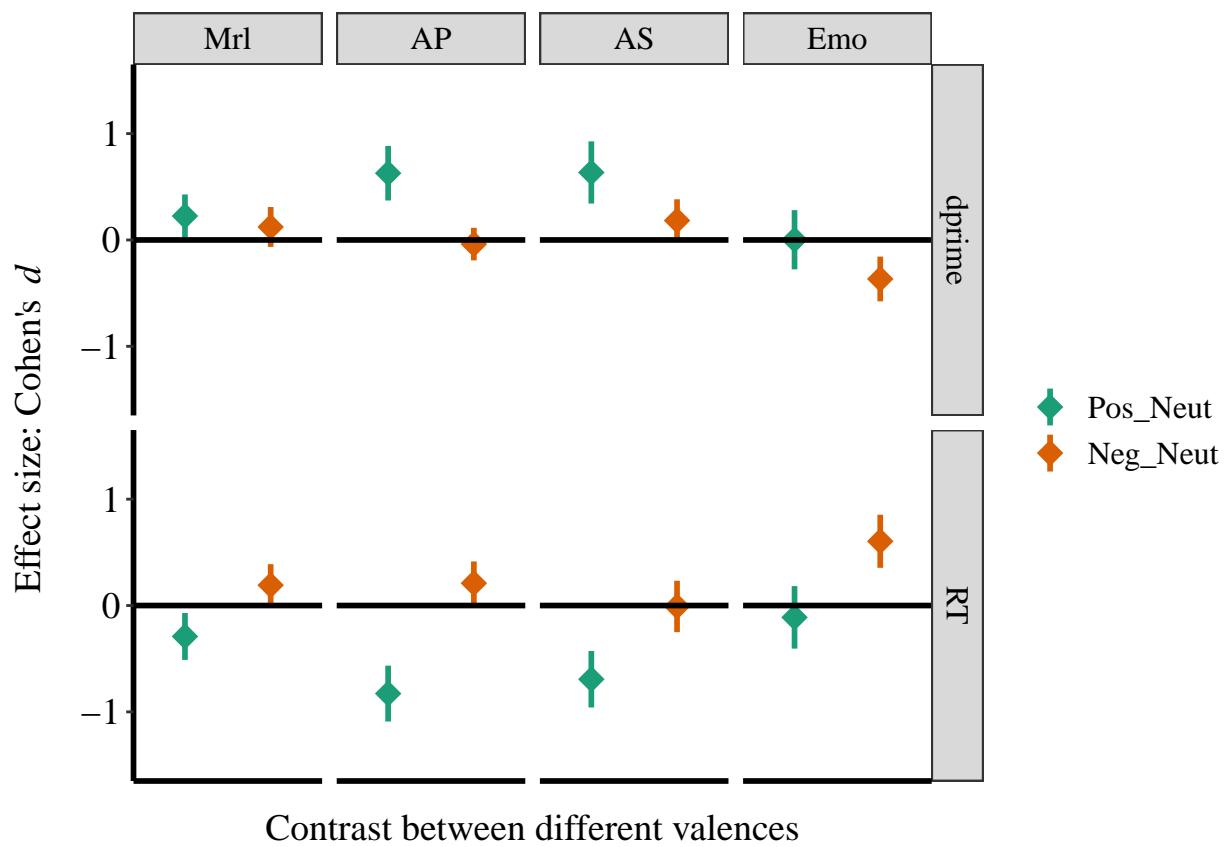


Figure 36. Effect size (Cohen's d) of Valence in Exp5.

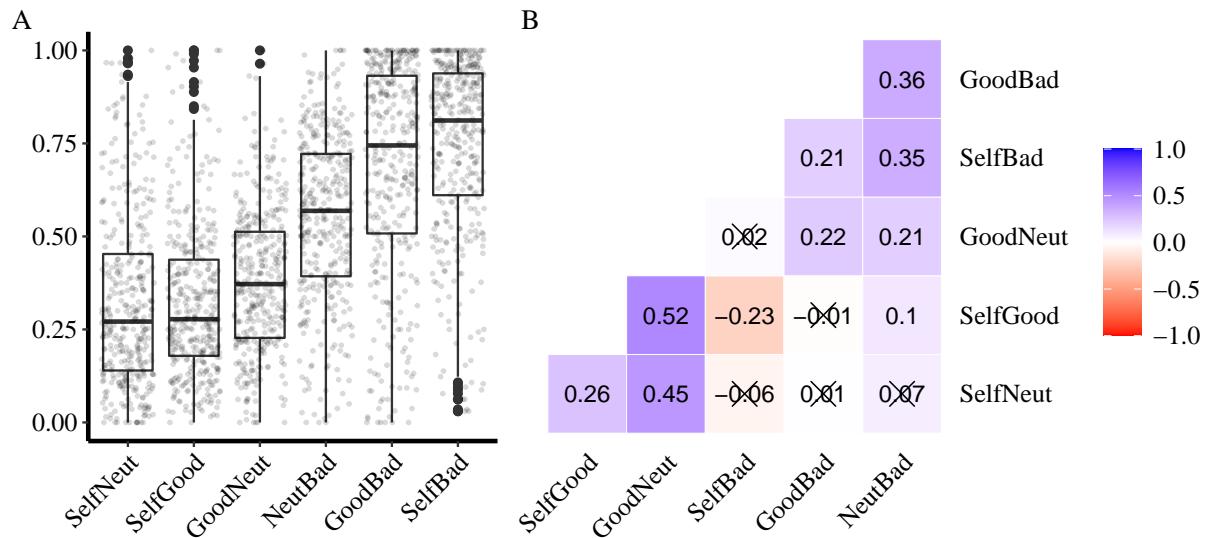


Figure 37. Self-rated personal distance

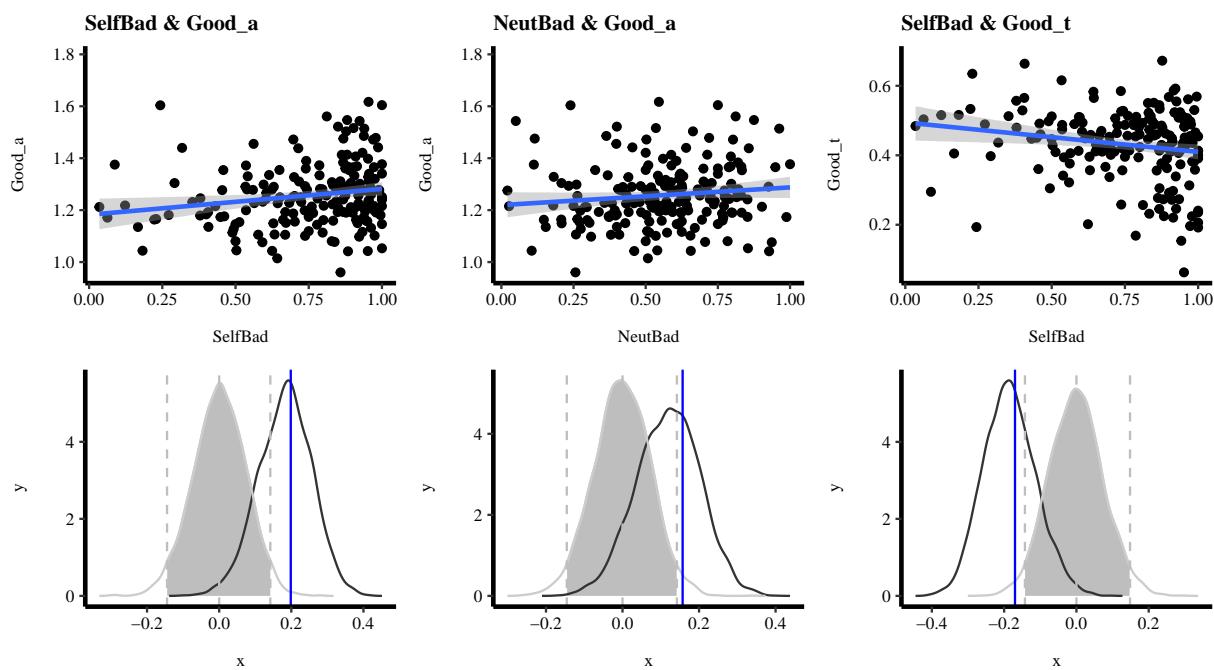


Figure 38. Correlation between moral identity and boundary separation of good condition; moral self-image and drift rate of good condition

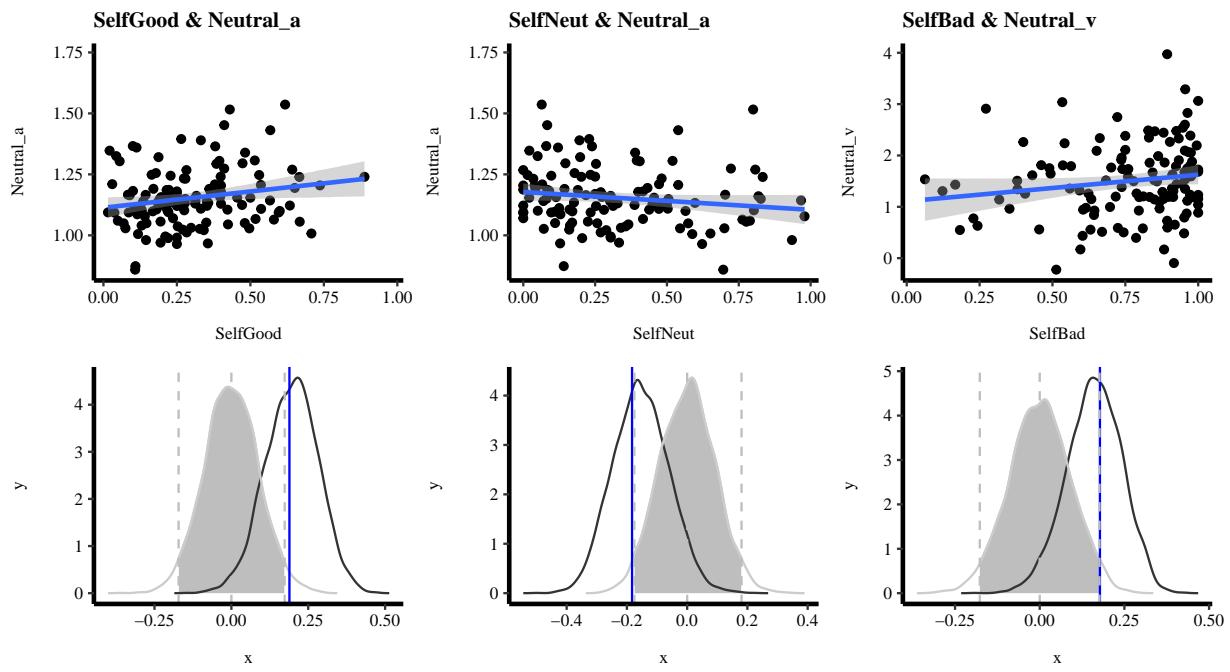


Figure 39. Correlation between personal distance and boundary separation of neutral condition