

¹ Good-self based social categorization in perceptual matching

² Hu Chuan-Peng^{1,2}, Kaiping Peng², & Jie Sui³

³ ¹ School of Psychology, Nanjing Normal University

⁴ ² Tsinghua University, 100084 Beijing, China

⁵ ³ University of Aberdeen, Aberdeen, Scotland

⁶ Author Note

⁷ Hu Chuan-Peng, School of Psychology, Nanjing Normal University, 210024 Nanjing,

⁸ China. Kaiping Peng, Department of Psychology, Tsinghua University, 100084 Beijing,

⁹ China. Jie Sui, School of Psychology, University of Aberdeen, Aberdeen, Scotland.

¹⁰ Authors contribution: HCP, JS, & KP design the study, HCP collected the data,

¹¹ HCP analyzed the data and drafted the manuscript. KP & JS supported this project.

¹² Correspondence concerning this article should be addressed to Hu Chuan-Peng,

¹³ School of Psychology, Nanjing Normal University, Ninghai Road 122, Gulou District,

¹⁴ 210024 Nanjing, China. E-mail: hcp4715@gmail.com

15

Abstract

16 To navigate in a complex social world, individuals are constantly evaluating others' moral
17 character. Also, they are managing a moral self-view that is aligned with their goals.
18 Though moral character in person perception and moral self-enhancement had been
19 extensively studied, the perceptual process of moral character is unkown, we examined
20 the influence of immediately acquired moral information on perceptual matching processing
21 by using social associative learning paradigm (self-tagging paradigm). In a series of
22 experiments, participants learned the concept of moral character and visual cues (shapes)
23 and then perform a perceptual matching task. The results showed that when geometric
24 shapes, without soical meaning, that associated with good moral character were prioritized.
25 This patterns of results were robust when we change different semantic words or using
26 behavioral history as an proxy of mroal character. Also, this patterns were robust in both
27 simutanously presentation and sequential presentation. We then examined two competing
28 explanation for this effect: value-based prioritization or social-categorization based
29 prioritization. We manipulated the identity of different moral character explicitly and
30 found that the good moral character effect was strong when for the self-referential
31 conditions but not for other-referential condition. We further tested the good-self based
32 social categorization by presenting the identity or moral character information as
33 task-irrelevant stimuli, so that we can distinguish between the unique good-self hypothesis
34 and a more general good-person based social categorization hypothesis. The evidence
35 suggested that human are more likely has a good-person based categorization instead of a
36 unique good-self. Finally, we explored whether the positivity effect only exist in moral
37 domain and found that this effect was not limited to moral domain but also aesthetic
38 domain, but not affective valence *per se*. Exploratory analyses on task-questionnaire
39 relationship found that there are weak correlation between self-bad distance and behavioral
40 pattern. These results suggest that there exist a social categorization in perceptual
41 decision-making, which is based on personal traits (moral character) but not affective

⁴² valence.

⁴³ *Keywords:* Perceptual decision-making, Self positivity bias, morality

⁴⁴ Word count: X

45 Good-self based social categorization in perceptual matching

46 **Introduction**

47 [sentences in bracket are key ideas]

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

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

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

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

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

90 Perceptual process of moral character

91 [exp1a, b, c, and exp2]

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

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

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

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

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

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

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

¹⁴⁹ **Morality as a social-categorization?**

¹⁵⁰ [possible explanations: person-based self-categorization vs. stimuli-based valence] The
¹⁵¹ robust pattern from our first four experiment suggested that there are some reliable
¹⁵² mechanisms underneath the effect. One possible explanation is the value-based attention,
¹⁵³ which suggested that valuable stimuli is prioritized in our low-level cognitive processes.
¹⁵⁴ Because positive moral character is potentially rewarding, e.g., potential cooperators, it is
¹⁵⁵ valuable to individuals and therefore being prioritized. There are also evidence consistent
¹⁵⁶ with this idea []. For example, XXX found that trustworthy faces attracted attention more
¹⁵⁷ than untrustworthy faces, probably because trustworthy faces are more likely to be the
¹⁵⁸ collaborative partners subsequent tasks, which will bring reward. This explanation has an
¹⁵⁹ implicit assumption, that is, participants were automatically viewing these stimuli as
¹⁶⁰ self-relevant (Juechems & Summerfield, 2019; Reicher & Hopkins, 2016) and
¹⁶¹ threatening/rewarding because of their semantic meaning. In this explanation, we will view
¹⁶² the moral concept, and the moral character represented by the concept, as objects and only
¹⁶³ judge whether they are rewarding/threatening or potentially rewarding/threatening to us.

¹⁶⁴ Another possibility is that we will perceive those moral character as person and
¹⁶⁵ automatic categorize whether they are ingroup or ougroup, that is, the social
¹⁶⁶ categorization process. This account assumed that moral character served as a way to
¹⁶⁷ categorize other. In the first four experiments' situation, the identity of the moral
¹⁶⁸ character is ambiguous, participants may automatically categorize morally good people as
¹⁶⁹ ingroup and therefore preferentially processed these information.

¹⁷⁰ However, the above four experiments can not distinguish between these two
¹⁷¹ possibilities, because the concept “good-peron” can both be rewarding and be categorized
¹⁷² as ingroup member, and previous studies using associative learning paradigm revealed that
¹⁷³ both rewarding stimuli (e.g., Sui, He, & Humphreys, 2012) and in-group information [F.
¹⁷⁴ Enock, Sui, Hewstone, and Humphreys (2018); enock_overlap_2020] are prioritized.

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

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

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

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

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

233 Key concepts and discussing points:

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

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

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

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

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

244 Morality as a way for social-categorization (McHugh, McGann, Igou, & Kinsella,
245 2019)? People are more likely to identify themselves with trustworthy faces (Verosky &
246 Todorov, 2010) (trustworthy faces has longer RTs).

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

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

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

253

Disclosures

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

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

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

269

General methods

270 **Design and Procedure**

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

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

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

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

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

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

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

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

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

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

352 Data analysis

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

³⁵⁶ then analyzed in two Bayesian approaches.

³⁵⁷ ***Bayesian hierarchical generalized linear model (BGLM).*** We first tested the
³⁵⁸ effect of experimental manipulation using Bayesian hierarchical generalized linear model
³⁵⁹ (BGLM), because it provided three advantages over the classic NHST approach (repeated
³⁶⁰ measure ANOVA or t-tests): first, Bayesian models use posterior distribution of parameter
³⁶¹ for statistical inference, therefore provided uncertainty in estimation (Rouder & Lu,
³⁶² 2005), second, BGLM can use distribution that fit the real distribution, which is the case
³⁶³ for reaction time data (Rousselet & Wilcox, 2019), third, BGLM also integrated different
³⁶⁴ levels of analysis, fully account the variability from each participants. We used the r
³⁶⁵ package **BRMs** (Bürkner, 2017) to build the model, which used Stan (Carpenter et al., 2017)
³⁶⁶ to sample from the posterior.

³⁶⁷ ***Signal detection theory.*** As in (Hu, Lan, Macrae, & Sui, 2020; Sui, He, &
³⁶⁸ Humphreys, 2012), we also used signal detection approach to analyze the accuracy data.
³⁶⁹ More specifically, we assume the match trials are signal and the non-match trials are noise.
³⁷⁰ To estimate the sensitivity and criterion of SDT, we adopted the Bayesian hierarchical
³⁷¹ GLM approach from (Rouder & Lu, 2005). When modelling the accuracy data for one
³⁷² participant, we assume that the accuracy of each trial is Bernoulli distributed (binomial
³⁷³ with 1 trial), with probability p_i that $y_i = 1$.

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

³⁷⁴ In the perceptual matching task, the probability p_i can then be modeled as a function of
³⁷⁵ the trial type:

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

³⁷⁶ The outcomes y_i are 0 if the participant responded “nonmatch” on trial i , 1 if they
³⁷⁷ responded “match.” The probability of the “match” response for trial i for a participant is

³⁷⁸ p_i . We then write the generalized linear model on the probits (z-scores; Φ , “Phi”) of ps . Φ
³⁷⁹ is the cumulative normal density function and maps z scores to probabilities. Given this
³⁸⁰ parameterization, the intercept of the model (β_0) is the standardized false alarm rate
³⁸¹ (probability of saying 1 when predictor is 0), which we take as our criterion c . The slope of
³⁸² the model (β_1) is the increase of saying 1 when predictor is 1, in z -scores, which is another
³⁸³ expression of d' . Therefore, $c = -zHR = -\beta_0$, and $d' = \beta_1$.

³⁸⁴ In each experiment, we had multiple participants, to estimate the group-level
³⁸⁵ parameters, we need to estimate parameters on individual level and the group level
³⁸⁶ parameter simultaneously. In this case, as above, we first assume that the outcome of each
³⁸⁷ trial is Bernoulli distributed, with probability p_{ij} that $y_{ij} = 1$.

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

³⁸⁸ And the the generalized linear model was re-written to include two levels:

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

³⁸⁹ The outcomes y_{ij} are 0 if participant j responded “nonmatch” on trial i , 1 if they
³⁹⁰ responded “match.” The probability of the “match” response for trial i for subject j is p_{ij} .
³⁹¹ We again can write the generalized linear model on the probits (z-scores; Φ , “Phi”) of ps .

³⁹² The subjective-specific intercepts ($\beta_0 = -zFAR$) and slopes ($\beta_1 = d'$) are describe
³⁹³ by multivariate normal with means and a covariance matrix for the parameters.

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

³⁹⁴ For experiments that had 2 (matching: match vs. non-match) by 3 (moral character:
³⁹⁵ good vs. neutral vs. bad), i.e., experiment 1a, 1b, 1c, 2, 5, and 6a, the formula for BGLM is
³⁹⁶ as follow:

```

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

```

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

```

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

```

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

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

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

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

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

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

$$\sigma_j \sim Cauchy()$$

412 Formula used for modeling the data as follow:

```

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

415     or

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

```

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

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

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

⁴³⁸ model all the experiment with similar design:

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

⁴³⁹ where

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

⁴⁴⁰ The outcomes y_{ijk} are 0 if participant j in experiment k responded “nonmatch” on trial i ,

⁴⁴¹ 1 if they responded “match.”

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

⁴⁴² and the experiment level parameter μ_{0k} and μ_{1k} is from a higher order

⁴⁴³ distribution:

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

⁴⁴⁴ in which μ_0 and μ_1 means the population level parameter.

⁴⁴⁵ In similar way, we expanded the RT model three-level model in which participants

⁴⁴⁶ and experiments are two group level variable and participants were nested in the

⁴⁴⁷ experiments.

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

⁴⁴⁸ y_{ijk} is the RT of the i th trial of the j th participants in the k th experiment.

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

⁴⁴⁹

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

450

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

451

$$\theta_k \sim Cauchy()$$

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

456 *Effect of moral character.* We synthesized effect size of d' and RT from experiment
 457 1a, 1b, 1c, 2, 5, and 6a for the effect of moral character. We reported the synthesized the
 458 effect across all experiments that tested the valence effect, using the mini meta-analysis
 459 approach (Goh, Hall, & Rosenthal, 2016).

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

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

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

473 ***Behavior-Questionnaire correlation.*** Finally, we explored correlation between
 474 results from behavioral results and self-reported measures.

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

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

486 Part 1: Perceptual processing moral character related inforation

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

490 Experiment 1a

491 Methods.

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

499 **Stimuli and Tasks.** Three geometric shapes were used in this experiment:

500 triangle, square, and circle. These shapes were paired with three labels (bad person, good
501 person or neutral person). The pairs were counterbalanced across participants.

502 **Procedure.** This experiment had two phases. First, there was a brief learning

503 stage. Participants were asked to learn the relationship between geometric shapes (triangle,
504 square, and circle) and different concepts of moral character (bad person, a good person, or
505 a neutral person). For example, a participant was told, “bad person is a circle; good person
506 is a triangle; and a neutral person is a square.” After participants remembered the
507 associations (usually in a few minutes), they started a practicing phase of matching task
508 which had the exact task as in the experimental task.

509 In the experimental task, participants judged whether shape–label pairs, which were

510 subsequently presented, were correct (i.e., the same as they learned). Each trial started
511 with the presentation of a central fixation cross for 500 ms. Subsequently, a pairing of a
512 shape and label (good person, bad person, and neutral person) was presented for 100 ms.

513 The pair presented could confirm to the verbal instruction for each pairing given in the
514 training stage, or it could be a recombination of a shape with a different label, with the
515 shape–label pairings being generated at random. The next frame showed a blank for
516 1100ms. Participants were expected to judge whether the shape was correctly assigned to
517 the person by pressing one of the two response buttons as quickly and accurately as
518 possible within this timeframe (to encourage immediate responding). Feedback (correct or
519 incorrect) was given on the screen for 500 ms at the end of each trial, if no response
520 detected, “too slow” was presented to remind participants to accelerate. Participants were
521 informed of their overall accuracy at the end of each block. The practice phase finished and
522 the experimental task began after the overall performance of accuracy during practice
523 phase achieved 60%.

524 For participants from the Tsinghua community, they completed 6 experimental blocks

525 of 60 trials. Thus, there were 60 trials in each condition (bad-person match, bad-person

526 nonmatch, good-person match, good-person nonmatch, neutral-person match, and
527 neutral-person nonmatch). For the participants from Wenzhou University, they finished 6
528 blocks of 120 trials, therefore, 120 trials for each condition.

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

533 **Results.**

534 **Hypothesis testing.**

535 *d prime.* We fitted a Bayesian hierarchical GLM for signal detection theory. The
536 results showed that when the shapes were tagged with labels with different moral character,
537 the sensitivity (d') and criteria (c) were both influenced. For the d' , we found that the
538 shapes associated with good person (2.46, 95% CI[2.21 2.72]) is greater than shapes tagged
539 with moral bad (2.07, 95% CI[1.83 2.32]), $P_{PosteriorComparison} = 1$. Shape tagged with
540 morally good person is also greater than shapes tagged with neutral person (2.23, 95%
541 CI[1.95 2.49]), $P_{PosteriorComparison} = 0.97$. Also, the shapes tagged with neutral person is
542 greater than shapes tagged with morally bad person, $P_{PosteriorComparison} = 0.92$.

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

548 **Reaction times.** We fitted a Bayesian hierarchical GLM for RTs, with a log-normal
549 distribution as the link function. We used the posterior distribution of the regression
550 coefficient to make statistical inferences. As in previous studies, the matched conditions are
551 much faster than the mismatched trials ($P_{PosteriorComparison} = 1$). We focused on matched

552 trials only, and compared different conditions: Good is faster than the neutral,
553 $P_{PosteriorComparison} = .99$, it was also faster than the Bad condition,
554 $P_{PosteriorComparison} = 1$. And the neutral condition is faster than the bad condition,
555 $P_{PosteriorComparison} = .99$. However, the mismatched trials are largely overlapped. See
556 Figure ??.

557 **HDDM.** We fitted our data with HDDM, using the response-coding (See also, Hu,
558 Lan, Macrae, & Sui, 2020). We estimated separate drift rate (v), non-decision time (T_0),
559 and boundary separation (a) for each condition. We found that the shapes tagged with
560 good person has higher drift rate and higher boundary separation than shapes tagged with
561 both neutral and bad person. Also, the shapes tagged with neutral person has a higher drift
562 rate than shapes tagged with bad person, but not for the boundary separation. Finally, we
563 found that shapes tagged with bad person had longer non-decision time (see Figure ??).

564 **Experiment 1b**

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

569 **Method.**

570 **Participants.** 72 college students (49 female, age = 20.17 ± 2.08 years)
571 participated. 39 of them were recruited from Tsinghua University community in 2014; 33
572 were recruited from Wenzhou University in 2017. All participants were right-handed except
573 one, and all had normal or corrected-to-normal vision. Informed consent was obtained from
574 all participants prior to the experiment according to procedures approved by the local
575 ethics committees. 20 participant's data were excluded from analysis because nearly
576 random level of accuracy, leaving 52 participants (36 female, age = 20.25 ± 2.31 years).

577 **Stimuli and Tasks.** Three geometric shapes (triangle, square, and circle, with 3.7°

578 $\times 3.7^\circ$ of visual angle) were presented above a white fixation cross subtending $0.8^\circ \times 0.8^\circ$

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

580 three labels with different moral valence: a morally bad person (“,” ERen), a morally

581 good person (“,” ShanRen) or a morally neutral person (“,” ChangRen). The order of

582 the associations between shapes and labels was counterbalanced across participants.

583 Three labels used in this experiment was selected based on the rating results from an

584 independent survey, in which participants rated the familiarity, frequency, and concreteness

585 of eight different words online. Of the eight words, three of them are morally positive

586 (HaoRen, ShanRen, Junzi), two of them are morally neutral (ChangRen, FanRen), and

587 three of them are morally negative (HuaiRen, ERen, LiuMang). An independent sample

588 consist of 35 participants (22 females, age 20.6 ± 3.11) were recruited to rate these words.

589 Based on the ratings (see supplementary materials Figure S1), we selected ShanRen,

590 ChangRen, and ERen to represent morally positive, neutral, and negative person.

591 **Procedure.** For participants from both Tsinghua community and Wenzhou

592 community, the procedure in the current study was exactly same as in experiment 1a.

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

594 **Results.**

595 **NHST.** Figure ?? shows d prime and reaction times of experiment 1b.

596 *d prime.* Repeated measures ANOVA revealed main effect of valence,

597 $F(1.83, 93.20) = 14.98$, $MSE = 0.18$, $p < .001$, $\hat{\eta}_G^2 = .053$. Paired t test showed that the

598 Good-Person condition (1.87 ± 0.102) was with greater d prime than Neutral condition

599 (1.44 ± 0.101 , $t(51) = 5.945$, $p < 0.001$). We also found that the Bad-Person condition

600 (1.67 ± 0.11) has also greater d prime than neutral condition , $t(51) = 3.132$, $p = 0.008$).

601 There Good-person condition was also slightly greater than the bad condition, $t(51) =$

602 2.265 , $p = 0.0701$.

603 *Reaction times.* We found interaction between Matchness and Valence

604 ($F(1.95, 99.31) = 19.71, MSE = 960.92, p < .001, \hat{\eta}_G^2 = .031$) and then analyzed the
 605 matched trials and mismatched trials separately, as in experiment 1a. For matched trials,
 606 we found the effect of valence $F(1.94, 99.10) = 33.97, MSE = 1,343.19, p < .001,$
 607 $\hat{\eta}_G^2 = .115$. Post-hoc *t*-tests revealed that shapes associated with Good Person (684 ± 8.77)
 608 were responded faster than Neutral-Person (740 ± 9.84), ($t(51) = -8.167, p < 0.001$) and
 609 Bad Person (728 ± 9.15), ($t(51) = -5.724, p < 0.0001$). While there was no significant
 610 differences between Neutral and Bad-Person condition ($t(51) = 1.686, p = 0.221$). For
 611 non-matched trials, there was no significant effect of Valence ($F(1.90, 97.13) = 1.80,$
 612 $MSE = 430.15, p = .173, \hat{\eta}_G^2 = .003$).

613 **BGLM.**

614 *Signal detection theory analysis of accuracy.* We fitted a Bayesian hierarchical GLM

615 for SDT. The results showed that when the shapes were tagged with labels with different
 616 moral valence, the sensitivity (d') and criteria (c) were both influence. For the d' , we found
 617 that the shapes tagged with morally good person (2.46, 95% CI[2.21 2.72]) is greater than
 618 shapes tagged with moral bad (2.07, 95% CI[1.83 2.32]), $P_{PosteriorComparison} = 1$. Shape
 619 tagged with morally good person is also greater than shapes tagged with neutral person
 620 (2.23, 95% CI[1.95 2.49]), $P_{PosteriorComparison} = 0.97$. Also, the shapes tagged with neutral
 621 person is greater than shapes tagged with morally bad person, $P_{PosteriorComparison} = 0.92$.

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

623 tagged with good person has the highest criteria (-1.01, [-1.14 -0.88]), followed by shapes
 624 tagged with neutral person(1.06, [-1.21 -0.92]), and then the shapes tagged with bad
 625 person(-1.11, [-1.25 -0.97]). However, pair-wise comparison showed that only showed strong
 626 evidence for the difference between good and bad conditions.

627 *Reaction time.* We fitted a Bayesian hierarchical GLM for RTs, with a log-normal

628 distribution as the link function. We used the posterior distribution of the regression

629 coefficient to make statistical inferences. As in previous studies, the matched conditions are
630 much faster than the mismatched trials ($P_{PosteriorComparison} = 1$). We focused on matched
631 trials only, and compared different conditions: Good is faster than the neutral,
632 $P_{PosteriorComparison} = .99$, it was also faster than the Bad condition,
633 $P_{PosteriorComparison} = 1$. And the neutral condition is faster than the bad condition,
634 $P_{PosteriorComparison} = .99$. However, the mismatched trials are largely overlapped. See
635 Figure ??.

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

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

646 Experiment 1c

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

651 Method.

652 **Participants.** 23 college students (15 female, age = 22.61 ± 2.62 years)
653 participated. All of them were recruited from Tsinghua University community in 2014.

654 Informed consent was obtained from all participants prior to the experiment according to
655 procedures approved by the local ethics committees. No participant was excluded because
656 they overall accuracy were above 0.6.

657 **Stimuli and Tasks.** Three geometric shapes (triangle, square, and circle, with
658 $3.7^\circ \times 3.7^\circ$ of visual angle) were presented above a white fixation cross subtending $0.8^\circ \times$
659 0.8° of visual angle at the center of the screen. The three most common names were
660 chosen, which are neutral in moral valence before the manipulation. Three names (Zhang,
661 Wang, Li) were first paired with three paragraphs of behavioral description. Each
662 description includes one sentence of biographic information and four sentences that
663 describing the moral behavioral under that name. To assess the that these three
664 descriptions represented good, neutral, and bad valence, we collected the ratings of three
665 person on six dimensions: morality, likability, trustworthiness, dominance, competence, and
666 aggressiveness, from an independent sample ($n = 34$, 18 female, age = 19.6 ± 2.05). The
667 rating results showed that the person with morally good behavioral description has higher
668 score on morality ($M = 3.59$, $SD = 0.66$) than neutral ($M = 0.88$, $SD = 1.1$), $t(33) =$
669 12.94 , $p < .001$, and bad conditions ($M = -3.4$, $SD = 1.1$), $t(33) = 30.78$, $p < .001$. Neutral
670 condition was also significant higher than bad conditions $t(33) = 13.9$, $p < .001$ (See
671 supplementary materials).

672 **Procedure.** After arriving the lab, participants were informed to complete two
673 experimental tasks, first a social memory task to remember three person and their
674 behaviors, after tested for their memory, they will finish a perceptual matching task. In the
675 social memory task, the descriptions of three person were presented without time
676 limitation. Participant self-paced to memorized the behaviors of each person. After they
677 memorizing, a recognition task was used to test their memory effect. Each participant was
678 required to have over 95% accuracy before preceding to matching task. The perceptual
679 learning task was followed, three names were randomly paired with geometric shapes.
680 Participants were required to learn the association and perform a practicing task before

681 they start the formal experimental blocks. They kept practicing until they reached 70%
 682 accuracy. Then, they would start the perceptual matching task as in experiment 1a. They
 683 finished 6 blocks of perceptual matching trials, each have 120 trials.

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

685 **Results.** Figure ?? shows d prime and reaction times of experiment 1c. We
 686 conducted same analysis as in Experiment 1a. Our analysis didn't show effect of valence
 687 on d prime, $F(1.93, 42.56) = 0.23$, $MSE = 0.41$, $p = .791$, $\hat{\eta}_G^2 = .005$. Neither the effect of
 688 valence on RT ($F(1.63, 35.81) = 0.22$, $MSE = 2,212.71$, $p = .761$, $\hat{\eta}_G^2 = .001$) or
 689 interaction between valence and matchness on RT ($F(1.79, 39.43) = 1.20$,
 690 $MSE = 1,973.91$, $p = .308$, $\hat{\eta}_G^2 = .005$).

691 **Signal detection theory analysis of accuracy.** We fitted a Bayesian
 692 hierarchical GLM for SDT. The results showed that when the shapes were tagged with
 693 labels with different moral valence, the sensitivity (d') and criteria (c) were both
 694 influenced. For the d' , we found that the shapes tagged with morally good person (2.30,
 695 95% CI[1.93 2.70]) is greater than shapes tagged with moral bad (2.11, 95% CI[1.83 2.42]),
 696 $P_{PosteriorComparison} = 0.8$. Shape tagged with morally good person is also greater than
 697 shapes tagged with neutral person (2.16, 95% CI[1.88 2.45]), $P_{PosteriorComparison} = 0.75$.

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

703 **Reaction time.** We fitted a Bayesian hierarchical GLM for RTs, with a log-normal
 704 distribution as the link function. We used the posterior distribution of the regression
 705 coefficient to make statistical inferences. As in previous studies, the matched conditions are
 706 much faster than the mismatched trials ($P_{PosteriorComparison} = .75$). We focused on

707 matched trials only, and compared different conditions: Good () is not faster than the
708 neutral (), $P_{PosteriorComparison} = .5$, it was faster than the Bad condition (),
709 $P_{PosteriorComparison} = .88$. And the neutral condition is faster than the bad condition,
710 $P_{PosteriorComparison} = .95$. However, the mismatched trials are largely overlapped.

711 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu,
712 Lan, Macrae, & Sui, 2020). We estimated separate drift rate (v), non-decision time (T_0),
713 and boundary separation (a) for each condition. We found that the shapes tagged with
714 good person has higher drift rate and higher boundary separation than shapes tagged with
715 both neutral and bad person. Also, the shapes tagged with neutral person has a higher drift
716 rate than shapes tagged with bad person, but not for the boundary separation. Finally, we
717 found that shapes tagged with bad person had longer non-decision time (see figure ??)).

718 Experiment 2: Sequential presenting

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

731 Method.

Participants. 35 participants (17 female, age = 21.66 ± 3.03) were recruited. 24

of them had participated in Experiment 1a (9 male, mean age = 21.9, s.d. = 2.9), and the

time gap between these experiment 1a and experiment 2 is at least six weeks. The results

of 1 participants were excluded from analysis because of less than 60% overall accuracy,

remains 34 participants (17 female, age = 21.74 ± 3.04).

Procedure. In Experiment 2, the sequential presenting makes the matching task

much easier than experiment 1. To avoid ceiling effect on behavioral data, we did a few

pilot experiments to get optimal parameters, i.e., the conditions under which participant

have similar accuracy as in Experiment 1 (around 70 ~ 80% accuracy). In the final

procedure, the label (good person, bad person, or neutral person) was presented for 50 ms

and then masked by a scrambled image for 200 ms. A geometric shape followed the

scrambled mask for 50 ms in a noisy background (which was produced by first

decomposing a square with $\frac{3}{4}$ gray area and $\frac{1}{4}$ white area to small squares with a size of 2

\times 2 pixels and then re-combine these small pieces randomly), instead of pure gray

background in Experiment 1. After that, a blank screen was presented 1100 ms, during

which participants should press a button to indicate the label and the shape match the

original association or not. Feedback was given, as in study 1. The next trial then started

after 700 ~ 1100 ms blank. Other aspects of study 2 were identical to study 1.

Data analysis. Data was analyzed as in study 1a.

Results.

NHST. Figure ?? shows d prime and reaction times of experiment 2. Less than

0.2% correct trials with less than 200ms reaction times were excluded.

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

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

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

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

758 no-significant difference between Neutral-person and Bad-person condition, $t(33) = -1.802$,
 759 $p = 0.185$.

760 *Reaction time.* The results of reaction times of matchness trials showed similar
 761 pattern as the d prime data.

762 We found interaction between Matchness and Valence ($F(1.99, 65.70) = 9.53$,
 763 $MSE = 605.36$, $p < .001$, $\hat{\eta}_G^2 = .017$) and then analyzed the matched trials and
 764 mismatched trials separately, as in experiment 1a. For matched trials, we found the effect
 765 of valence $F(1.99, 65.76) = 10.57$, $MSE = 1,192.65$, $p < .001$, $\hat{\eta}_G^2 = .067$. Post-hoc t -tests
 766 revealed that shapes associated with Good Person (548 ± 9.4) were responded faster than
 767 Neutral-Person (582 ± 10.9), ($t(33) = -3.95$, $p = 0.0011$) and Bad Person (582 ± 10.2),
 768 $t(33) = -3.9$, $p = 0.0013$). While there was no significant differences between Neutral and
 769 Bad-Person condition ($t(33) = -0.01$, $p = 0.999$). For non-matched trials, there was no
 770 significant effect of Valence ($F(1.99, 65.83) = 0.17$, $MSE = 489.80$, $p = .843$, $\hat{\eta}_G^2 = .001$).

771 **BGLMM.**

772 *Signal detection theory analysis of accuracy.* We fitted a Bayesian hierarchical GLM
 773 for SDT. The results showed that when the shapes were tagged with labels with different
 774 moral valence, the sensitivity (d') and criteria (c) were both influence. For the d' , we found
 775 that the shapes tagged with morally good person (2.46, 95% CI[2.21 2.72]) is greater than
 776 shapes tagged with moral bad (2.07, 95% CI[1.83 2.32]), $P_{PosteriorComparison} = 1$. Shape
 777 tagged with morally good person is also greater than shapes tagged with neutral person
 778 (2.23, 95% CI[1.95 2.49]), $P_{PosteriorComparison} = 0.97$. Also, the shapes tagged with neutral
 779 person is greater than shapes tagged with morally bad person, $P_{PosteriorComparison} = 0.92$.

780 Interesting, we also found the criteria for three conditions also differ, the shapes
 781 tagged with good person has the highest criteria (-1.01, [-1.14 -0.88]), followed by shapes
 782 tagged with neutral person(1.06, [-1.21 -0.92]), and then the shapes tagged with bad
 783 person(-1.11, [-1.25 -0.97]). However, pair-wise comparison showed that only showed strong

784 evidence for the difference between good and bad conditions.

785 *Reaction times.* We fitted a Bayesian hierarchical GLM for RTs, with a log-normal
786 distribution as the link function. We used the posterior distribution of the regression
787 coefficient to make statistical inferences. As in previous studies, the matched conditions are
788 much faster than the mismatched trials ($P_{PosteriorComparison} = .75$). We focused on
789 matched trials only, and compared different conditions: Good () is not faster than the
790 neutral (), $P_{PosteriorComparison} = .5$, it was faster than the Bad condition (),
791 $P_{PosteriorComparison} = .88$. And the neutral condition is faster than the bad condition,
792 $P_{PosteriorComparison} = .95$. However, the mismatched trials are largely overlapped.

793 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu,
794 Lan, Macrae, & Sui, 2020). We estimated separate drift rate (v), non-decision time (T_0),
795 and boundary separation (a) for each condition. We found that the shapes tagged with
796 good person has higher drift rate and higher boundary separation than shapes tagged with
797 both neutral and bad person. Also, the shapes tagged with neutral person has a higher
798 drift rate than shapes tagged with bad person, but not for the boundary separation.
799 Finally, we found that shapes tagged with bad person had longer non-decision time (see
800 figure @ref(fig:plot-exp1c -HDDM))).

801 Discussion

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

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

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

811 **Method.**

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

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

820 **Stimuli.** Three geometric shapes (triangle, square and circle, each $4.6^\circ \times 4.6^\circ$ of
821 visual angle) were presented at the center of screen for 50 ms after 500ms of fixation (0.8°
822 $\times 0.8^\circ$ of visual angle). The association of the three shapes to bad person (“ , HuaiRen”),
823 good person (“ , HaoRen”) or ordinary person (“ , ChangRen”) was counterbalanced
824 across participants. The words bad person, good person or ordinary person ($3.6^\circ \times 1.6^\circ$)
825 was also displayed at the center fo the screen. Participants had to judge whether the
826 pairings of label and shape matched (e.g., Does the circle represent a bad person?). The
827 experiment was run on a PC using E-prime software (version 2.0). These stimuli were
828 displayed on a 22-in CRT monitor (1024×768 at 100Hz). We used backward masking to
829 avoid over-processing of the moral words, in which a scrambled picture were presented for
830 900 ms after the label. Also, to avoid the ceiling effect on accuracy, shapes were presented
831 on a noisy background based on our pilot studies. The noisy images were made by
832 scrambling a picture of 3/4gray and 1/4 white at resolution of 2×2 pixel.

833 **Procedure.** The procedure was similar to Experiment 2. Participants finished 9

834 blocks of trial, each with 120 trials. In total, participants finished 180 trials for each

835 combination of condition.

836 As in experiment 2 (Sui, He, & Humphreys, 2012), subjects first learned the

837 associations between labels and shapes and then completed a shape-label matching task

838 (e.g., good person-triangle). In each trial of the matching task, a fixation were first

839 presented for 500 ms, followed by a 50 ms label; then, a scrambled picture presented 900

840 ms. After the backward mask, the shape were presented on a noisy background for 50ms.

841 Participant have to response in 1000ms after the presentation of the shape, and finally, a

842 feedback screen was presented for 500 ms (see figure 1). The inter-trial interval (ITI) were

843 randomly varied at the range of 1000 ~ 1400 ms.

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

845 2.0 was used to present stimuli and collect behavioral results. Data were collected and

846 analyzed when accuracy performance in total reached 60%.

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

848 **Results.**

849 **NHST.** Only the behavioral results were reported here. Figure ?? shows *d* prime

850 and reaction times of experiment 6a.

851 *d* prime. We conducted repeated measures ANOVA, with moral valence as

852 independent variable. The results revealed the main effect of valence

853 ($F(1.74, 40.05) = 3.76$, $MSE = 0.10$, $p = .037$, $\eta^2_G = .021$). Post-hoc analysis revealed that

854 shapes link with Good person (mean = 3.13, SE = 0.109) is greater than Neutral condition

855 (mean = 2.88, SE = 0.14), $t = 2.916$, $df = 24$, $p = 0.02$, p-value adjusted by Tukey method,

856 but the *d* prime between Good and bad (mean = 3.03, SE = 0.142) ($t = 1.512$, $df = 24$, $p =$

857 = 0.3034, p-value adjusted by Tukey method), bad and neutral ($t = 1.599$, $df = 24$, $p =$

858 0.2655, p-value adjusted by Tukey method) were not significant.

859 *Reaction times.* The results of reaction times of matchness trials showed similar

860 pattern as the d prime data.

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

862 $MSE = 450.47, p < .001, \hat{\eta}_G^2 = .021$) and then analyzed the matched trials and

863 mismatched trials separately, as in experiment 2. For matched trials, we found the effect of

864 valence $F(1.97, 45.25) = 32.37, MSE = 522.42, p < .001, \hat{\eta}_G^2 = .078$. For non-matched

865 trials, there was no significant effect of Valence ($F(1.77, 40.67) = 0.35, MSE = 242.15,$

866 $p = .679, \hat{\eta}_G^2 = .000$). Post-hoc t -tests revealed that shapes associated with Good Person

867 (mean = 550, SE = 13.8) were responded faster than Neutral-Person (501, SE = 14.7),

868 ($t(24) = -5.171, p = 0.0001$) and Bad Person (523, SE = 16.3), $t(24) = -8.137, p <$

869 0.0001.), and Neutral is faster than Bad-Person condition ($t(32) = -3.282, p = 0.0085$).

870 **BGLM.**

871 *Signal detection theory analysis of accuracy.*

872 *Reaction time.*

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

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

875 and boundary separation (a) for each condition. We found that, similar to experiment 2,

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

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

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

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

880 for the self-referential condition.

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

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

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

884 non-match trials (see figure ??).

Part 2: interaction between valence and identity

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

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

893 Method.

894 **Participants.** 38 college students (15 female, age = 21.92 ± 2.16) participated in
895 experiment 3a. All of them were right-handed, and all had normal or corrected-to-normal
896 vision. Informed consent was obtained from all participants prior to the experiment
897 according to procedures approved by a local ethics committee. One female and one male
898 student did not finish the experiment, and 1 participants' data were excluded from analysis
899 because less than 60% overall accuracy, remains 35 participants (13 female, age = $22.11 \pm$
900 2.13).

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

906 **Stimuli.** The stimuli used in study 3a share the same parameters with experiment
907 1 & 2. The differences was that we used six shapes: triangle, square, circle, trapezoid,
908 diamond, regular pentagon, and six labels: good self, neutral self, bad self, good person,
909 bad person, and neutral person. To match the concreteness of the label, we asked
910 participant to chosen an unfamiliar name of their own gender to be the stranger.

911 **Procedure.** After being fully explained and signed the informed consent,
912 participants were instructed to chose a name that can represent a stranger with same
913 gender as the participant themselves, from a common Chinese name pool. Before
914 experiment, the experimenter explained the meaning of each label to participants. For
915 example, the “good self” mean the morally good side of themselves, them could imagine
916 the moment when they do something’s morally applauded, “bad self” means the morally
917 bad side of themselves, they could also imagine the moment when they doing something
918 morally wrong, and “neutral self” means the aspect of self that does not related to
919 morality, they could imagine the moment when they doing something irrelevant to
920 morality. In the same sense, the “good other,” “bad other,” and “neutral other” means the
921 three different aspects of the stranger, whose name was chosen before the experiment.
922 Then, the experiment proceeded as study 1a. Each participant finished 6 blocks, each have
923 120 trials. The sequence of trials was pseudo-randomized so that there are 10 matched
924 trials for each condition and 10 non-matched trials for each condition (good self, neutral
925 self, bad self, good other, neutral other, bad other) for each block.

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

929 Results.

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

932 d prime. There was evidence for the main effect of valence, $F(1.89, 64.37) = 11.09$,
933 $MSE = 0.23$, $p < .001$, $\hat{\eta}_G^2 = .039$, and main effect of self-relevance, $F(1, 34) = 3.22$,
934 $MSE = 0.54$, $p = .082$, $\hat{\eta}_G^2 = .015$, as well as the interaction, $F(1.79, 60.79) = 3.39$,
935 $MSE = 0.43$, $p = .045$, $\hat{\eta}_G^2 = .022$.

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

937 The valence effect was shown for the self-referential conditions, $F(1.65, 56.25) = 13.98$,
 938 $MSE = 0.31$, $p < .001$, $\hat{\eta}_G^2 = .119$. Post-hoc test revealed that the Good-Self condition
 939 (1.97 ± 0.14) was with greater d prime than Neutral condition (1.41 ± 0.12 , $t(34) = 4.505$,
 940 $p = 0.0002$), and Bad-self condition (1.43 ± 0.102), $t(34) = 3.856$, $p = 0.0014$. There was
 941 difference between neutral and bad condition, $t(34) = -0.238$, $p = 0.9694$. However, no
 942 effect of valence was found for the other-referential condition $F(1.98, 67.36) = 0.38$,
 943 $MSE = 0.35$, $p = .681$, $\hat{\eta}_G^2 = .004$.

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

947 For the match trials, we found that the interaction between identity and valence,
 948 $F(1.72, 58.61) = 3.89$, $MSE = 2,750.19$, $p = .032$, $\hat{\eta}_G^2 = .019$, as well as the main effect of
 949 valence $F(1.98, 67.34) = 35.76$, $MSE = 1,127.25$, $p < .001$, $\hat{\eta}_G^2 = .079$, but not the effect of
 950 identity $F(1, 34) = 0.20$, $MSE = 3,507.14$, $p = .660$, $\hat{\eta}_G^2 = .001$. As for the d prime, we
 951 separated analyzed the self-referential and other-referential trials. For the Self-referential
 952 trials, we found the main effect of valence, $F(1.80, 61.09) = 30.39$, $MSE = 1,584.53$,
 953 $p < .001$, $\hat{\eta}_G^2 = .159$; for the other-referential trials, the effect of valence is weaker,
 954 $F(1.86, 63.08) = 2.85$, $MSE = 2,224.30$, $p = .069$, $\hat{\eta}_G^2 = .024$. We then focused on the self
 955 conditions: the good-self condition (713 ± 12) is faster than neutral- (776 ± 11.8), $t(34) =$
 956 -7.396 , $p < .0001$, and bad-self (772 ± 10.1) conditions, $t(34) = -5.66$, $p < .0001$. But
 957 there is not difference between neutral- and bad-self conditions, $t(34) = 0.481$, $p = 0.881$.

958 For the nonmatch trials, we didn't found any strong effect: identity, $F(1, 34) = 3.43$,
 959 $MSE = 660.02$, $p = .073$, $\hat{\eta}_G^2 = .004$, valence $F(1.89, 64.33) = 0.40$, $MSE = 444.10$,
 960 $p = .661$, $\hat{\eta}_G^2 = .001$, or interaction between the two $F(1.94, 66.02) = 2.42$, $MSE = 817.35$,
 961 $p = .099$, $\hat{\eta}_G^2 = .007$.

962 **BGLM.**

963 *Signal detection theory analysis of accuracy.* We found that the d prime is greater
964 when shapes were associated with good self condition than with neutral self or bad self, but
965 shapes associated with bad self and neutral self didn't show differences. Comparing the self
966 vs other under three condition revealed that shapes associated with good self is greater
967 than with good other, but with a weak evidence. In contrast, for both neutral and bad
968 valence condition, shapes associated with other had greater d prime than with self.

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

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

980 **Experiment 3b**

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

986 **Participants.** Study 3b were finished in 2017, at that time we have calculated that
987 the effect size (Cohen's d) of good-person (or good-self) vs. bad-person (or bad-other) was

988 between 0.47 ~ 0.53, based on the data from Tsinghua community in study 1a, 1b, 2, 3a,
989 4a, and 4b. Based on this effect size, we estimated that 54 participants would allow we to
990 detect the effect size of Cohen's = 0.5 with 95% power and alpha = 0.05, using G*power
991 3.192 (Faul, Erdfelder, Buchner, & Lang, 2009). Therefore, we planned to stop after we
992 arrived this number. During the data collected at Wenzhou University, 61 participants (45
993 females; 19 to 25 years of age, age = 20.42 ± 1.77) came to the testing room and we tested
994 all of them during a single day. All participants were right-handed, and all had normal or
995 corrected-to-normal vision. Informed consent was obtained from all participants prior to
996 the experiment according to procedures approved by a local ethics committee. 4
997 participants' data were excluded from analysis because their over all accuracy was lower
998 than 60%, 1 more participant was excluded because of zero hit rate for one condition,
999 leaving 56 participants (43 females; 19 to 25 years old, age = 20.27 ± 1.60).

1000 **Design.** Study 3b has the same experimental design as 3a, with a $2 \times 3 \times 2$
1001 within-subject design. The first variable was self-relevance, include two levels: self-relevant
1002 vs. stranger-relevant; the second variable was moral valence, include good, neutral and bad;
1003 the third variable was the matching between shape and label: match vs. mismatch.
1004 Stimuli. The stimuli used in study 3b share the same parameters with experiment 3a. 6
1005 shapes were included (triangle, square, circle, trapezoid, diamond, regular pentagon), as
1006 well as 6 labels, but the labels changed to "good self," "neutral self," "bad self," "good
1007 him/her," bad him/her", "neutral him/her," the stranger's label is consistent with
1008 participants' gender. Same as study 3a, we asked participant to chosen an unfamiliar name
1009 of their own gender to be the stranger before showing them the relationship. Note, because
1010 of implementing error, the personal distance data did not collect for this experiment.

1011 **Stimuli.** The stimuli used in study 3b is the same as in experiment 3a.

1012 **Procedure.** In this experiment, participants finished two matching tasks, i.e.,
1013 self-matching task, and other-matching task. In the self-matching task, participants first
1014 associate the three aspects of self to three different shapes, and then perform the matching

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

1023 Data Analysis. Same as experiment 3a.

1024 Results.

1025 NHST. Figure 6 shows d prime and reaction times of experiment 3b. Less than 5%
1026 correct trials with less than 200ms reaction times were excluded.

1027 d prime. There was no evidence for the main effect of valence,
1028 $F(1.92, 105.43) = 1.90$, $MSE = 0.33$, $p = .157$, $\hat{\eta}_G^2 = .005$, but we found a main effect of
1029 self-relevance, $F(1, 55) = 4.65$, $MSE = 0.89$, $p = .035$, $\hat{\eta}_G^2 = .017$, as well as the interaction,
1030 $F(1.90, 104.36) = 5.58$, $MSE = 0.26$, $p = .006$, $\hat{\eta}_G^2 = .011$.

1031 We then conducted separated ANOVA for self-referential and other-referential trials.
1032 The valence effect was shown for the self-referential conditions, $F(1.75, 96.42) = 6.73$,
1033 $MSE = 0.30$, $p = .003$, $\hat{\eta}_G^2 = .037$. Post-hoc test revealed that the Good-Self condition
1034 (2.15 ± 0.12) was with greater d prime than Neutral condition (1.83 ± 0.12 , $t(34) = 3.36$,
1035 $p = 0.0031$), and Bad-self condition (1.87 ± 0.12), $t(34) = 2.955$, $p = 0.01$. There was
1036 difference between neutral and bad condition, $t(34) = -0.039$, $p = 0.914$. However, no effect
1037 of valence was found for the other-referential condition $F(1.93, 105.97) = 0.61$,
1038 $MSE = 0.31$, $p = .539$, $\hat{\eta}_G^2 = .002$.

1039 Reaction time. We found interaction between Matchness and Valence
1040 ($F(1.86, 102.47) = 15.44$, $MSE = 3, 112.78$, $p < .001$, $\hat{\eta}_G^2 = .006$) and then analyzed the

¹⁰⁴¹ matched trials and nonmatch trials separately, as in previous experiments.

¹⁰⁴² For the match trials, we found that the interaction between identity and valence,
¹⁰⁴³ $F(1.67, 92.11) = 6.14$, $MSE = 6,472.48$, $p = .005$, $\hat{\eta}_G^2 = .009$, as well as the main effect of
¹⁰⁴⁴ valence $F(1.88, 103.65) = 24.25$, $MSE = 5,994.25$, $p < .001$, $\hat{\eta}_G^2 = .038$, but not the effect
¹⁰⁴⁵ of identity $F(1, 55) = 48.49$, $MSE = 25,892.59$, $p < .001$, $\hat{\eta}_G^2 = .153$. As for the d prime,
¹⁰⁴⁶ we separated analyzed the self-referential and other-referential trials. For the
¹⁰⁴⁷ Self-referential trials, we found the main effect of valence, $F(1.66, 91.38) = 23.98$,
¹⁰⁴⁸ $MSE = 6,965.61$, $p < .001$, $\hat{\eta}_G^2 = .100$; for the other-referential trials, the effect of valence
¹⁰⁴⁹ is weaker, $F(1.89, 103.94) = 5.96$, $MSE = 5,589.90$, $p = .004$, $\hat{\eta}_G^2 = .014$. We then focused
¹⁰⁵⁰ on the self conditions: the good-self condition (713 ± 12) is faster than neutral- ($776 \pm$
¹⁰⁵¹ 11.8), $t(34) = -7.396$, $p < .0001$, and bad-self (772 ± 10.1) conditions, $t(34) = -5.66$, $p <$
¹⁰⁵² $.0001$. But there is not difference between neutral- and bad-self conditions, $t(34) = 0.481$, p
¹⁰⁵³ $= 0.881$.

¹⁰⁵⁴ For the nonmatch trials, we didn't found any strong effect: identity, $F(1, 55) = 10.31$,
¹⁰⁵⁵ $MSE = 24,590.52$, $p = .002$, $\hat{\eta}_G^2 = .035$, valence $F(1.98, 108.63) = 20.57$, $MSE = 2,847.51$,
¹⁰⁵⁶ $p < .001$, $\hat{\eta}_G^2 = .016$, or interaction between the two $F(1.93, 106.25) = 35.51$,
¹⁰⁵⁷ $MSE = 1,939.88$, $p < .001$, $\hat{\eta}_G^2 = .019$.

¹⁰⁵⁸ **BGLM.**

¹⁰⁵⁹ *Signal detection theory analysis of accuracy.* We found that the d prime is greater
¹⁰⁶⁰ when shapes were associated with good self condition than with neutral self or bad self, but
¹⁰⁶¹ shapes associated with bad self and neutral self didn't show differences. comparing the self
¹⁰⁶² vs other under three condition revealed that shapes associated with good self is greater
¹⁰⁶³ than with good other, but with a weak evidence. In contrast, for both neutral and bad
¹⁰⁶⁴ valence condition, shapes associated with other had greater d prime than with self.

¹⁰⁶⁵ *Reaction time.* In reaction times, we found that same trends in the match trials as
¹⁰⁶⁶ in the RT: while the shapes associated with good self was greater than with good other

1067 (log mean diff = -0.02858, 95%HPD[-0.070898, 0.0154]), the direction is reversed for
1068 neutral and negative condition. see Figure 7

1069 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu,
1070 Lan, Macrae, & Sui, 2020). We estimated separate drift rate (v), non-decision time (T_0),
1071 and boundary separation (a) for each condition. We found that, similar to experiment 3a,
1072 the shapes tagged with good person has higher drift rate and higher boundary separation
1073 than shapes tagged with both neutral and bad person, but only for the self-referential
1074 condition. Also, the shapes tagged with neutral person has a higher drift rate than shapes
1075 tagged with bad person, but not for the boundary separation, and this effect also exist only
1076 for the self-referential condition.

1077 Interestingly, we found that in both self-referential and other-referential conditions,
1078 the shapes associated bad valence have higher drift rate and higher boundary separation.
1079 which might suggest that the shape associated with bad stimuli might be prioritized in the
1080 non-match trials (see figure 8)).

1081 **Experiment 6b**

1082 Experiment 6b was conducted to study the neural correlates of the prioritization
1083 effect of positive self, i.e., the neural underlying of the behavioral effect found int
1084 experiment 3a. However, as in experiment 6a, the procedure of this experiment was
1085 modified to adopted to ERP experiment.

1086 **Method.**

1087 **Participants.** 23 college students (8 female, age = 22.86 ± 2.47) participated the
1088 current study, all of them were recruited from Tsinghua University in 2016. Informed
1089 consent was obtained from all participants prior to the experiment according to procedures
1090 approved by a local ethics committee. For day 1's data, 1 participant was excluded from
1091 the current analysis because of lower than 60% overall accuracy, remaining 22 participants

1092 (8 female, age = 22.76 ± 2.49). For day 2's data, one participant dropped out, leaving 22
1093 participants (9 female, age = 23.05 ± 2.46), all of them has overall accuracy higher than
1094 60%.

1095 **Design.** The experimental design of this experiment is same as experiment 3: a 2
1096 $\times 3 \times 2$ within-subject design with self-relevance (self-relevant vs. other-relevant), moral
1097 valence (good, neutral, and bad) and matchness between shape and label (match
1098 vs. mismatch) as within-subject variables.

1099 **Stimuli.** As in experiment 3a, 6 shapes were included (triangle, square, circle,
1100 trapezoid, diamond, regular pentagon), as well as 6 labels (good self, neutral self, bad self,
1101 good person, bad person, neutral person). To match the concreteness of the label, we asked
1102 participant to chosen an unfamiliar name of their own gender to be the stranger.

1103 **Procedure.** The procedure was similar to Experiment 2 and 6a. Subjects first
1104 learned the associations between labels and shapes and then completed a shape-label
1105 matching task. In each trial of the matching task, a fixation were first presented for 500
1106 ms, followed by a 50 ms label; then, a scrambled picture presented 900 ms. After the
1107 backward mask, the shape were presented on a noisy background for 50ms. Participant
1108 have to response in 1000ms after the presentation of the shape, and finally, a feedback
1109 screen was presented for 500 ms. The inter-trial interval (ITI) were randomly varied at the
1110 range of 1000 ~ 1400 ms.

1111 All the stimuli were presented on a gray background (RGB: 127, 127, 127). E-primed
1112 2.0 was used to present stimuli and collect behavioral results. Data were collected and
1113 analyzed when accuracy performance in total reached 60%.

1114 Because learning 6 associations was more difficult than 3 associations and participant
1115 might have low accuracy (see experiment 3a), the current study had extended to a two-day
1116 paradigm to maximizing the accurate trials that can be used in EEG data. At the first day,
1117 participants learnt the associations and finished 9 blocks of the matching task, each had

₁₁₁₈ 120 trials, without EEG recording. That is, each condition has 90 trials.

₁₁₁₉ Participants came back to lab at the second day and finish the same task again, with
₁₁₂₀ EEG recorded. Before the EEG experiment, each participant finished a practice session
₁₁₂₁ again, if their accuracy is equal or higher than 85%, they start the experiment (one
₁₁₂₂ participant used lower threshold 75%). Each participant finished 18 blocks, each has 90
₁₁₂₃ trials. One participant finished additional 6 blocks because of high error rate at the
₁₁₂₄ beginning, another two participant finished addition 3 blocks because of the technique
₁₁₂₅ failure in recording the EEG data. To increase the number of trials that can be used for
₁₁₂₆ EEG data analysis, matched trials has twice number as mismatched trials, therefore, for
₁₁₂₇ matched trials each participants finished 180 trials for each condition, for mismatched
₁₁₂₈ trials, each conditions has 90 trials.

₁₁₂₉ **Data Analysis.** Same as experiment 3a.

₁₁₃₀ **Results of Day 1.**

₁₁₃₁ **NHST.** Figure 9 shows d prime and reaction times of experiment 3b. Less than 5%
₁₁₃₂ correct trials with less than 200ms reaction times were excluded.

₁₁₃₃ *d prime.* There was no evidence for the main effect of valence,
₁₁₃₄ $F(1.91, 40.20) = 11.98$, $MSE = 0.15$, $p < .001$, $\hat{\eta}_G^2 = .040$, but we found a main effect of
₁₁₃₅ self-relevance, $F(1, 21) = 1.21$, $MSE = 0.20$, $p = .284$, $\hat{\eta}_G^2 = .003$, as well as the interaction,
₁₁₃₆ $F(1.28, 26.90) = 12.88$, $MSE = 0.21$, $p = .001$, $\hat{\eta}_G^2 = .041$.

₁₁₃₇ We then conducted separated ANOVA for self-referential and other-referential trials.
₁₁₃₈ The valence effect was shown for the self-referential conditions, $F(1.73, 36.42) = 29.31$,
₁₁₃₉ $MSE = 0.14$, $p < .001$, $\hat{\eta}_G^2 = .147$. Post-hoc test revealed that the Good-Self condition
₁₁₄₀ (2.15 ± 0.12) was with greater d prime than Neutral condition (1.83 ± 0.12 , $t(34) = 3.36$,
₁₁₄₁ $p = 0.0031$), and Bad-self condition (1.87 ± 0.12), $t(34) = 2.955$, $p = 0.01$. There was
₁₁₄₂ difference between neutral and bad condition, $t(34) = -0.039$, $p = 0.914$. However, no effect
₁₁₄₃ of valence was found for the other-referential condition $F(1.75, 36.72) = 0.00$, $MSE = 0.18$,

₁₁₄₄ $p = .999$, $\hat{\eta}_G^2 = .000$.

₁₁₄₅ *Reaction time.* We found interaction between Matchness and Valence

₁₁₄₆ ($F(1.79, 37.63) = 4.07$, $MSE = 704.90$, $p = .029$, $\hat{\eta}_G^2 = .003$) and then analyzed the

₁₁₄₇ matched trials and nonmatch trials separately, as in previous experiments.

₁₁₄₈ For the match trials, we found that the interaction between identity and valence,

₁₁₄₉ $F(1.72, 36.16) = 4.55$, $MSE = 1,560.90$, $p = .022$, $\hat{\eta}_G^2 = .015$, as well as the main effect of

₁₁₅₀ valence $F(1.93, 40.55) = 9.83$, $MSE = 1,951.84$, $p < .001$, $\hat{\eta}_G^2 = .044$, but not the effect of

₁₁₅₁ identity $F(1, 21) = 4.87$, $MSE = 2,032.05$, $p = .039$, $\hat{\eta}_G^2 = .012$. As for the d prime, we

₁₁₅₂ separated analyzed the self-referential and other-referential trials. For the Self-referential

₁₁₅₃ trials, we found the main effect of valence, $F(1.92, 40.38) = 14.48$, $MSE = 1,647.20$,

₁₁₅₄ $p < .001$, $\hat{\eta}_G^2 = .112$; for the other-referential trials, the effect of valence is weaker,

₁₁₅₅ $F(1.79, 37.50) = 1.04$, $MSE = 1,842.07$, $p = .356$, $\hat{\eta}_G^2 = .008$. We then focused on the self

₁₁₅₆ conditions: the good-self condition (713 ± 12) is faster than neutral- (776 ± 11.8), $t(34) =$

₁₁₅₇ -7.396 , $p < .0001$, and bad-self (772 ± 10.1) conditions, $t(34) = -5.66$, $p < .0001$. But

₁₁₅₈ there is not difference between neutral- and bad-self conditions, $t(34) = 0.481$, $p = 0.881$.

₁₁₅₉ For the nonmatch trials, we didn't found any strong effect: identity, $F(1, 21) = 2.76$,

₁₁₆₀ $MSE = 1,718.93$, $p = .112$, $\hat{\eta}_G^2 = .006$, valence $F(1.61, 33.77) = 3.81$, $MSE = 1,532.21$,

₁₁₆₁ $p = .041$, $\hat{\eta}_G^2 = .012$, or interaction between the two $F(1.90, 39.97) = 2.23$, $MSE = 720.80$,

₁₁₆₂ $p = .123$, $\hat{\eta}_G^2 = .004$.

₁₁₆₃ **BGLM.**

₁₁₆₄ *Signal detection theory analysis of accuracy.* We found that the d prime is greater

₁₁₆₅ when shapes were associated with good self condition than with neutral self or bad self, but

₁₁₆₆ shapes associated with bad self and neutral self didn't show differences. comparing the self

₁₁₆₇ vs other under three condition revealed that shapes associated with good self is greater

₁₁₆₈ than with good other, but with a weak evidence. In contrast, for both neutral and bad

₁₁₆₉ valence condition, shapes associated with other had greater d prime than with self.

1170 *Reaction time.* In reaction times, we found that same trends in the match trials as
1171 in the RT: while the shapes associated with good self was greater than with good other
1172 ($\log \text{mean diff} = -0.02858$, 95%HPD[-0.070898, 0.0154]), the direction is reversed for
1173 neutral and negative condition. see Figure 10

1174 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu,
1175 Lan, Macrae, & Sui, 2020). We estimated separate drift rate (v), non-decision time (T_0),
1176 and boundary separation (a) for each condition. We found that, similar to experiment 3a,
1177 the shapes tagged with good person has higher drift rate and higher boundary separation
1178 than shapes tagged with both neutral and bad person, but only for the self-referential
1179 condition. Also, the shapes tagged with neutral person has a higher drift rate than shapes
1180 tagged with bad person, but not for the boundary separation, and this effect also exist only
1181 for the self-referential condition.

1182 Interestingly, we found that in both self-referential and other-referential conditions,
1183 the shapes associated bad valence have higher drift rate and higher boundary separation.
1184 which might suggest that the shape associated with bad stimuli might be prioritized in the
1185 non-match trials (see figure 11).

1186 **Part 3: Implicit binding between valence and identity**

1187 In this part, we reported two studies in which the moral valence or the self-referential
1188 processing is not task-relevant. We are interested in testing whether the task-relevance will
1189 eliminate the effect observed in previous experiment.

1190 **Experiment 4a: Morality as task-irrelevant variable**

1191 In part two (experiment 3a and 3b), participants learned the association between self
1192 and moral valence directly. In Experiment 4a, we examined whether the interaction
1193 between moral valence and identity occur even when one of the variable was irrelevant to

the task. In experiment 4a, participants learnt associations between shapes and self/other labels, then made perceptual match judgments only about the self or other conditions labels and shapes (cf. Sui, He, and Humphreys (2012)). However, we presented labels of different moral valence in the shapes, which means that the moral valence factor become task irrelevant. If the binding between moral good and self is intrinsic and automatic, then we will observe that facilitating effect of moral good for self conditions, but not for other conditions.

Method.

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

Design. As in Experiment 3, a $2 \times 3 \times 2$ within-subject design was used. The first variable was self-relevance (self and stranger associations); the second variable was moral valence (good, neutral and bad associations); the third variable was the matching between shape and label (matching vs. non-match for the personal association). However, in this 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,

1221 neutral-other, and bad-other.

1222 ***Stimuli.*** 2 shapes were included (circle, square) and each appeared above a central
 1223 fixation cross with the personal label appearing below. However, the shapes were not
 1224 empty but with a two-Chinese-character word in the middle, the word was one of three
 1225 labels with different moral valence: “good person,” “bad person” and “neutral person.”
 1226 Before the experiment, participants learned the self/other association, and were informed
 1227 to only response to the association between shapes’ configures and the labels below the
 1228 fixation, but ignore the words within shapes. Besides the behavioral experiments,
 1229 participants from Tsinghua community also finished questionnaires as Experiments 3, and
 1230 participants from Wenzhou community finished a series of questionnaire as the other
 1231 experiment finished in Wenzhou.

1232 ***Procedure.*** The procedure was similar to Experiment 1. There were 6 blocks of
 1233 trial, each with 120 trials for 2017 data. Due to procedure error, the data collected in 2015
 1234 in Tsinghua community only have 60 trials for each block, i.e., 30 trials per condition.

1235 As in study 3a, before each task, the instruction showed the meaning of each label to
 1236 participants. The self-matching task and other-matching task were randomized between
 1237 participants. Each participant finished 6 blocks, each have 120 trials.

1238 ***Data Analysis.*** Same as experiment 3a.

1239 **Results.**

1240 ***NHST.*** Figure 12 shows d prime and reaction times of experiment 3a. Less than
 1241 5% correct trials with less than 200ms reaction times were excluded.

1242 ***d prime.*** There was no evidence for the main effect of valence,
 1243 $F(1.93, 111.66) = 0.53$, $MSE = 0.12$, $p = .581$, $\hat{\eta}_G^2 = .000$, but we found a main effect of
 1244 self-relevance, $F(1, 58) = 121.04$, $MSE = 0.48$, $p < .001$, $\hat{\eta}_G^2 = .189$, as well as the
 1245 interaction, $F(1.99, 115.20) = 4.12$, $MSE = 0.14$, $p = .019$, $\hat{\eta}_G^2 = .004$.

1246 We then conducted separated ANOVA for self-referential and other-referential trials.

1247 The valence effect was shown for the self-referential conditions, $F(1.95, 112.92) = 3.01$,
 1248 $MSE = 0.15$, $p = .055$, $\hat{\eta}_G^2 = .008$. Post-hoc test revealed that the Good-Self condition
 1249 (2.15 ± 0.12) was with greater d prime than Neutral condition (1.83 ± 0.12 , $t(34) = 3.36$,
 1250 $p = 0.0031$), and Bad-self condition (1.87 ± 0.12), $t(34) = 2.955$, $p = 0.01$. There was
 1251 difference between neutral and bad condition, $t(34) = -0.039$, $p = 0.914$. However, no effect
 1252 of valence was found for the other-referential condition $F(1.98, 114.61) = 1.75$,
 1253 $MSE = 0.10$, $p = .179$, $\hat{\eta}_G^2 = .003$.

1254 *Reaction time.* We found interaction between Matchness and Valence

1255 ($F(1.94, 112.64) = 0.84$, $MSE = 465.35$, $p = .432$, $\hat{\eta}_G^2 = .000$) and then analyzed the
 1256 matched trials and nonmatch trials separately, as in previous experiments.

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

1258 $F(1.90, 110.18) = 4.41$, $MSE = 465.91$, $p = .016$, $\hat{\eta}_G^2 = .003$, as well as the main effect of
 1259 valence $F(1.98, 114.82) = 0.94$, $MSE = 606.30$, $p = .392$, $\hat{\eta}_G^2 = .001$, but not the effect of
 1260 identity $F(1, 58) = 124.15$, $MSE = 4,037.53$, $p < .001$, $\hat{\eta}_G^2 = .257$. As for the d prime, we
 1261 separated analyzed the self-referential and other-referential trials. For the Self-referential
 1262 trials, we found the main effect of valence, $F(1.97, 114.32) = 6.29$, $MSE = 367.25$,
 1263 $p = .003$, $\hat{\eta}_G^2 = .006$; for the other-referential trials, the effect of valence is weaker,
 1264 $F(1.95, 112.89) = 0.35$, $MSE = 699.50$, $p = .699$, $\hat{\eta}_G^2 = .001$. We then focused on the self
 1265 conditions: the good-self condition (713 ± 12) is faster than neutral- (776 ± 11.8), $t(34) =$
 1266 -7.396 , $p < .0001$, and bad-self (772 ± 10.1) conditions, $t(34) = -5.66$, $p < .0001$. But
 1267 there is not difference between neutral- and bad-self conditions, $t(34) = 0.481$, $p = 0.881$.

1268 For the nonmatch trials, we didn't found any strong effect: identity, $F(1, 58) = 0.16$,

1269 $MSE = 1,547.37$, $p = .692$, $\hat{\eta}_G^2 = .000$, valence $F(1.96, 113.52) = 0.68$, $MSE = 390.26$,
 1270 $p = .508$, $\hat{\eta}_G^2 = .000$, or interaction between the two $F(1.90, 110.27) = 0.04$,
 1271 $MSE = 585.80$, $p = .953$, $\hat{\eta}_G^2 = .000$.

BGLM.

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

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

HDDM. We fitted our data with HDDM, using the response-coding (also see Hu, Lan, Macrae, & Sui, 2020). 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 14)).

Experiment 4b: Morality as task-irrelevant variable

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

1297 performance of good self condition, even identity is irrelevant to the task.

1298 **Method.**

1299 **Participants.** 53 participants (39 female, age = 20.57 ± 1.81) participated the
1300 current study, 34 of them were from Tsinghua University in 2015, 19 were from Wenzhou
1301 University participated in 2017. All participants were right-handed, and all had normal or
1302 corrected-to-normal vision. Informed consent was obtained from all participants prior to
1303 the experiment according to procedures approved by a local ethics committee. The data
1304 from 8 participants from Wenzhou site were excluded from analysis because their accuracy
1305 was close to chance (< 0.6). The results for the remaining 45 participants (33 female, age
1306 = 20.78 ± 1.76) were analyzed and reported.

1307 **Design.** As in Experiment 3, a $2 \times 3 \times 2$ within-subject design was used. The first
1308 variable was self-relevance (self and stranger associations); the second variable was moral
1309 valence (good, neutral and bad associations); the third variable was the matching between
1310 shape and label (matching vs. non-match for the personal association). However, in this
1311 the task, participants only learn the association between two geometric shapes and two
1312 labels (self and other), i.e., only self-relevance were related to the task. The moral valence
1313 manipulation was achieved by embedding the personal label of the labels in the geometric
1314 shapes, see below. For simplicity, the trials where shapes where paired with self and with a
1315 word of “good person” inside were shorted as good-self condition, similarly, the trials where
1316 shapes paired with the self and with a word of “bad person” inside were shorted as bad-self
1317 condition. Hence, we also have six conditions: good-self, neutral-self, bad-self, good-other,
1318 neutral-other, and bad-other.

1319 **Stimuli.** 2 shapes were included (circle, square) and each appeared above a central
1320 fixation cross with the personal label appearing below. However, the shapes were not
1321 empty but with a two-Chinese-character word in the middle, the word was one of three
1322 labels with different moral valence: “good person,” “bad person” and “neutral person.”

1323 Before the experiment, participants learned the self/other association, and were informed
 1324 to only response to the association between shapes' configures and the labels below the
 1325 fixation, but ignore the words within shapes. Besides the behavioral experiments,
 1326 participants from Tsinghua community also finished questionnaires as Experiments 3, and
 1327 participants from Wenzhou community finished a series of questionnaire as the other
 1328 experiment finished in Wenzhou.

1329 ***Procedure.*** The procedure was similar to Experiment 1. There were 6 blocks of
 1330 trial, each with 120 trials for 2017 data. Due to procedure error, the data collected in 2015
 1331 in Tsinghua community only have 60 trials for each block, i.e., 30 trials per condition.

1332 As in study 3a, before each task, the instruction showed the meaning of each label to
 1333 participants. The self-matching task and other-matching task were randomized between
 1334 participants. Each participant finished 6 blocks, each have 120 trials.

1335 ***Data Analysis.*** Same as experiment 3a.

1336 **Results.**

1337 ***NHST.*** Figure 15 shows d prime and reaction times of experiment 3a. Less than
 1338 5% correct trials with less than 200ms reaction times were excluded.

1339 d prime. There was no evidence for the main effect of valence,
 1340 $F(1.59, 69.94) = 2.34$, $MSE = 0.48$, $p = .115$, $\hat{\eta}_G^2 = .010$, but we found a main effect of
 1341 self-relevance, $F(1, 44) = 0.00$, $MSE = 0.08$, $p = .994$, $\hat{\eta}_G^2 = .000$, as well as the interaction,
 1342 $F(1.96, 86.41) = 0.53$, $MSE = 0.10$, $p = .585$, $\hat{\eta}_G^2 = .001$.

1343 We then conducted separated ANOVA for self-referential and other-referential trials.
 1344 The valence effect was shown for the self-referential conditions, $F(1.75, 76.86) = 3.08$,
 1345 $MSE = 0.25$, $p = .058$, $\hat{\eta}_G^2 = .017$. Post-hoc test revealed that the Good-Self condition
 1346 (2.15 ± 0.12) was with greater d prime than Neutral condition (1.83 ± 0.12 , $t(34) = 3.36$,
 1347 $p = 0.0031$), and Bad-self condition (1.87 ± 0.12), $t(34) = 2.955$, $p = 0.01$. There was
 1348 difference between neutral and bad condition, $t(34) = -0.039$, $p = 0.914$. However, no effect

₁₃₄₉ of valence was found for the other-referential condition $F(1.63, 71.50) = 1.07$, $MSE = 0.33$,
₁₃₅₀ $p = .336$, $\hat{\eta}_G^2 = .006$.

₁₃₅₁ *Reaction time.* We found interaction between Matchness and Valence
₁₃₅₂ ($F(1.87, 82.50) = 18.58$, $MSE = 1,291.12$, $p < .001$, $\hat{\eta}_G^2 = .023$) and then analyzed the
₁₃₅₃ matched trials and nonmatch trials separately, as in previous experiments.

₁₃₅₄ For the match trials, we found that the interaction between identity and valence,
₁₃₅₅ $F(1.86, 81.84) = 5.22$, $MSE = 308.30$, $p = .009$, $\hat{\eta}_G^2 = .003$, as well as the main effect of
₁₃₅₆ valence $F(1.80, 79.37) = 11.04$, $MSE = 2,937.54$, $p < .001$, $\hat{\eta}_G^2 = .059$, but not the effect of
₁₃₅₇ identity $F(1, 44) = 0.23$, $MSE = 263.26$, $p = .632$, $\hat{\eta}_G^2 = .000$. As for the d prime, we
₁₃₅₈ separated analyzed the self-referential and other-referential trials. For the Self-referential
₁₃₅₉ trials, we found the main effect of valence, $F(1.74, 76.48) = 13.69$, $MSE = 1,732.08$,
₁₃₆₀ $p < .001$, $\hat{\eta}_G^2 = .079$; for the other-referential trials, the effect of valence is weaker,
₁₃₆₁ $F(1.87, 82.44) = 7.09$, $MSE = 1,527.43$, $p = .002$, $\hat{\eta}_G^2 = .043$. We then focused on the self
₁₃₆₂ conditions: the good-self condition (713 ± 12) is faster than neutral- (776 ± 11.8), $t(34) =$
₁₃₆₃ -7.396 , $p < .0001$, and bad-self (772 ± 10.1) conditions, $t(34) = -5.66$, $p < .0001$. But
₁₃₆₄ there is not difference between neutral- and bad-self conditions, $t(34) = 0.481$, $p = 0.881$.

₁₃₆₅ For the nonmatch trials, we didn't found any strong effect: identity, $F(1, 44) = 1.96$,
₁₃₆₆ $MSE = 319.47$, $p = .169$, $\hat{\eta}_G^2 = .001$, valence $F(1.69, 74.54) = 6.59$, $MSE = 886.19$,
₁₃₆₇ $p = .004$, $\hat{\eta}_G^2 = .010$, or interaction between the two $F(1.88, 82.57) = 0.31$, $MSE = 316.96$,
₁₃₆₈ $p = .718$, $\hat{\eta}_G^2 = .000$.

₁₃₆₉ **BGLM.**

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

¹³⁷⁵ valence condition, shapes associated with other had greater d prime than with self.

¹³⁷⁶ *Reaction time.* In reaction times, we found that same trends in the match trials as
¹³⁷⁷ in the RT: while the shapes associated with good self was greater than with good other
¹³⁷⁸ ($\log \text{mean diff} = -0.02858$, 95%HPD[-0.070898, 0.0154]), the direction is reversed for
¹³⁷⁹ neutral and negative condition. see Figure 16

¹³⁸⁰ **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu,
¹³⁸¹ Lan, Macrae, & Sui, 2020). We estimated separate drift rate (v), non-decision time (T_0),
¹³⁸² and boundary separation (a) for each condition. We found that the shapes tagged with
¹³⁸³ good person has higher drift rate and higher boundary separation than shapes tagged with
¹³⁸⁴ both neutral and bad person. Also, the shapes tagged with neutral person has a higher drift
¹³⁸⁵ rate than shapes tagged with bad person, but not for the boundary separation. Finally, we
¹³⁸⁶ found that shapes tagged with bad person had longer non-decision time (see figure 17)).

¹³⁸⁷ Results

¹³⁸⁸ Effect of moral valence

¹³⁸⁹ In this part, we synthesized results from experiment 1a, 1b, 1c, 2, 5 and 6a. Data
¹³⁹⁰ from 192 participants were included in these analyses. We found differences between
¹³⁹¹ positive and negative conditions on RT was Cohen's $d = -0.58 \pm 0.06$, 95% CI [-0.70 -0.47];
¹³⁹² on d' was Cohen's $d = 0.24 \pm 0.05$, 95% CI [0.15 0.34]. The effect was also observed
¹³⁹³ between positive and neutral condition, RT: Cohen's $d = -0.44 \pm 0.10$, 95% CI [-0.63
¹³⁹⁴ -0.25]; d' : Cohen's $d = 0.31 \pm 0.07$, 95% CI [0.16 0.45]. And the difference between neutral
¹³⁹⁵ and bad conditions are not significant, RT: Cohen's $d = 0.15 \pm 0.07$, 95% CI [0.00 0.30];
¹³⁹⁶ d' : Cohen's $d = 0.07 \pm 0.07$, 95% CI [-0.08 0.21]. See Figure 18 left panel.

1397 Interaction between valence and self-reference

1398 In this part, we combined the experiments that explicitly manipulated the
1399 self-reference and valence, which includes 3a, 3b, 6b, 7a, and 7b. For the positive versus
1400 negative contrast, data were from five experiments with 178 participants; for positive
1401 versus neutral and neutral versus negative contrasts, data were from three experiments (1402 3a, 3b, and 6b) with 108 participants.

1403 In most of these experiments, the interaction between self-reference and valence was
1404 significant (see results of each experiment in supplementary materials). In the
1405 mini-meta-analysis, we analyzed the valence effect for self-referential condition and
1406 other-referential condition separately.

1407 For the self-referential condition, we found the same pattern as in the first part of
1408 results. That is we found significant differences between positive and neutral as well as
1409 positive and negative, but not neutral and negative. The effect size of RT between positive
1410 and negative is Cohen's $d = -0.89 \pm 0.12$, 95% CI [-1.11 -0.66]; on d' was Cohen's $d = 0.61$
1411 ± 0.09 , 95% CI [0.44 0.78]. The effect was also observed between positive and neutral
1412 condition, RT: Cohen's $d = -0.76 \pm 0.13$, 95% CI [-1.01 -0.50]; d' : Cohen's $d = 0.69 \pm$
1413 0.14 , 95% CI [0.42 0.96]. And the difference between neutral and bad conditions are not
1414 significant, RT: Cohen's $d = 0.03 \pm 0.13$, 95% CI [-0.22 0.29]; d' : Cohen's $d = 0.08 \pm 0.08$,
1415 95% CI [-0.07 0.24]. See Figure 18 the middle panel.

1416 For the other-referential condition, we found that only the difference between positive
1417 and negative on RT was significant, all the other conditions were not. The effect size of RT
1418 between positive and negative is Cohen's $d = -0.28 \pm 0.05$, 95% CI [-0.38 -0.17]; on d' was
1419 Cohen's $d = -0.02 \pm 0.08$, 95% CI [-0.17 0.13]. The effect was not observed between
1420 positive and neutral condition, RT: Cohen's $d = -0.12 \pm 0.10$, 95% CI [-0.31 0.06]; d' :
1421 Cohen's $d = 0.01 \pm 0.08$, 95% CI [-0.16 0.17]. And the difference between neutral and bad
1422 conditions are not significant, RT: Cohen's $d = 0.14 \pm 0.09$, 95% CI [-0.03 0.31]; d' :

¹⁴²³ Cohen's $d = 0.05 \pm 0.07$, 95% CI [-0.08 0.18]. See Figure 18 right panel.

¹⁴²⁴ **Generalizability of the valence effect**

¹⁴²⁵ In this part, we reported the results from experiment 4 in which either moral valence
¹⁴²⁶ or self-reference were manipulated as task-irrelevant stimuli.

¹⁴²⁷ For experiment 4a, when self-reference was the target and moral valence was
¹⁴²⁸ task-irrelevant, we found that only under the implicit self-referential condition, i.e., when
¹⁴²⁹ the moral words were presented as task irrelevant stimuli, there was the main effect of
¹⁴³⁰ valence and interaction between valence and reference for both d prime and RT (See
¹⁴³¹ supplementary results for the detailed statistics). For d prime, we found good-self
¹⁴³² condition (2.55 ± 0.86) had higher d prime than bad-self condition (2.38 ± 0.80); good self
¹⁴³³ condition was also higher than neutral self (2.45 ± 0.78) but there was not statistically
¹⁴³⁴ significant, while the neutral-self condition was higher than bad self condition and not
¹⁴³⁵ significant neither. For reaction times, good-self condition (654.26 ± 67.09) were faster
¹⁴³⁶ relative to bad-self condition (665.64 ± 64.59), and over neutral-self condition ($664.26 \pm$
¹⁴³⁷ 64.71). The difference between neutral-self and bad-self conditions were not significant.
¹⁴³⁸ However, for the other-referential condition, there was no significant differences between
¹⁴³⁹ different valence conditions. See Figure 19.

¹⁴⁴⁰ For experiment 4b, when valence was the target and the identity was task-irrelevant,
¹⁴⁴¹ we found a strong valence effect (see supplementary results and Figure 20, Figure 21).

¹⁴⁴² In this experiment, the advantage of good-self condition can only be disentangled by
¹⁴⁴³ comparing the self-referential and other-referential conditions. Therefore, we calculated the
¹⁴⁴⁴ differences between the valence effect under self-referential and other referential conditions
¹⁴⁴⁵ and used the weighted variance as the variance of this differences. We found this
¹⁴⁴⁶ modulation effect on RT. The valence effect of RT was stronger in self-referential than
¹⁴⁴⁷ other-referential for the Good vs. Neutral condition (-0.33 ± 0.01), and to a less extent the

¹⁴⁴⁸ Good vs. Bad condition (-0.17 ± 0.01). While the size of the other effect's CI included
¹⁴⁴⁹ zero, suggesting those effects didn't differ from zero. See Figure 22.

¹⁴⁵⁰ **Specificity of valence effect**

¹⁴⁵¹ In this part, we analyzed the results from experiment 5, which included positive,
¹⁴⁵² neutral, and negative valence from four different domains: morality, emotion, aesthetics of
¹⁴⁵³ human, and aesthetics of scene. We found interaction between valence and domain for both
¹⁴⁵⁴ *d* prime and RT (match trials). A common pattern appeared in all four domains: each
¹⁴⁵⁵ domain showed a binary results instead of gradient on both *d* prime and RT. For morality,
¹⁴⁵⁶ aesthetics of human, and aesthetics of scene, there was a positivity effect where the positive
¹⁴⁵⁷ conditions had advantages over both neutral (greater *d* prime and faster RT), while neutral
¹⁴⁵⁸ and negative conditions didn't differ from each other. But for the emotional stimuli, there
¹⁴⁵⁹ was a reversed negativity effect: positive and neutral conditions were not significantly
¹⁴⁶⁰ different from each other but both had advantage over negative conditions. See
¹⁴⁶¹ supplementary materials for detailed statistics. Also note that the effect size in moral
¹⁴⁶² domain is smaller than the aesthetic domains (beauty of people and beauty of scene). See
¹⁴⁶³ Figure 23.

¹⁴⁶⁴ **Self-reported personal distance**

¹⁴⁶⁵ See Figure 24.

¹⁴⁶⁶ **Correlation analyses**

¹⁴⁶⁷ The reliability of questionnaires can be found in (Liu et al., 2020). We calculated the
¹⁴⁶⁸ correlation between the data from behavioral task and the questionnaire data. First, we
¹⁴⁶⁹ calculated the score for each scale based on their structure and factor loading, instead of

1470 sum score (McNeish & Wolf, 2020). Then, we used SEM to estimate the correlation
1471 because it can include measurement model and statistical model in a unified framework.

1472 To make sure that what we found were not false positive, we used two method to
1473 ensure the robustness of our analysis. first, we split the data into two half: the data with
1474 self and without, then, we used the conditional random forest to find the robust correlation
1475 in the exploratory data (with self reference) that can be replicated in the confirmatory data
1476 (without the self reference). The robust correlation were then analyzed using SEM

1477 Instead of use the exploratory correlation analysis, we used a more principled way to
1478 explore the correlation between parameter of HDDM (v , t , and a) and scale scores and
1479 person distance.

1480 We didn't find the correlation between scale scores and the parameters of HDDM,
1481 but found weak correlation between personal distance and the parameter estimated from
1482 Good and neutral conditions.

1483 First, boundary separation (a) of moral good condition was correlated with both
1484 Self-Bad distance ($r = 0.198$, 95% CI $[], p = 0.0063$) and Neutral-Bad distance
1485 ($r = 0.1571$, 95% CI $[], p = 0.031$). At the same time, the non-decision time is negatively
1486 correlated with Self-Bad distance ($r = 0.169$, 95% CI $[], p = 0.0197$). See Figure 25.

1487 Second, we found the boundary separation of neutral condition is positively
1488 correlated with the personal distance between self and good distance ($r = 0.189$, 95% CI $[],$
1489 $p = 0.036$), but negatively correlated with self-neutral distance($r = -0.183$, 95% CI $[],$
1490 $p = 0.042$). Also, the drift rate of the neutral condition is positively correlated with the
1491 Self-Bad distance ($r = 0.177$, 95% CI $[], p = 0.048$).a. See figure 26

1492 We also explored the correlation between behavioral data and questionnaire scores
1493 separately for experiments with and without self-referential, however, the sample size is
1494 very low for some conditions.

1495

Discussion

1496

References

1497

Abele, A. E., Ellemers, N., Fiske, S. T., Koch, A., & Yzerbyt, V. (2020). Navigating the social world: Toward an integrated framework for evaluating self, individuals, and groups. *Psychological Review*. <https://doi.org/10.1037/rev0000262>

1500

Amodio, D. M. (2019). Social cognition 2.0: An interactive memory systems account. *Trends in Cognitive Sciences*, 23(1), 21–33.
<https://doi.org/10.1016/j.tics.2018.10.002>

1503

Anderson, E., Siegel, E. H., Bliss-Moreau, E., & Barrett, L. F. (2011). The visual impact of gossip. *Science*, 332(6036), 1446–1448.
<https://doi.org/10.1126/science.1201574>

1506

Brainard, D. H. (1997). The psychophysics toolbox [Journal Article]. *Spatial Vision*, 10(4), 433–436.

1508

Bürkner, P.-C. (2017). Brms: An r package for bayesian multilevel models using stan [Journal Article]. *Journal of Statistical Software; Vol 1, Issue 1 (2017)*. Retrieved from
<https://www.jstatsoft.org/v080/i01%20http://dx.doi.org/10.18637/jss.v080.i01>

1512

Carlson, R. W., Maréchal, M. A., Oud, B., Fehr, E., & Crockett, M. J. (2020). Motivated misremembering of selfish decisions. *Nature Communications*, 11(1), 2100. <https://doi.org/10.1038/s41467-020-15602-4>

1515

Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M., ... Riddell, A. (2017). Stan: A probabilistic programming language [Journal Article]. *Journal of Statistical Software*, 76(1).
<https://doi.org/10.18637/jss.v076.i01>

- 1519 Dang, J., King, K. M., & Inzlicht, M. (2020). Why are self-report and behavioral
1520 measures weakly correlated? *Trends in Cognitive Sciences*, 24(4), 267–269.
1521 <https://doi.org/10.1016/j.tics.2020.01.007>
- 1522 Dzhelyova, M., Perrett, D. I., & Jentzsch, I. (2012). Temporal dynamics of
1523 trustworthiness perception. *Brain Research*, 1435, 81–90.
1524 <https://doi.org/10.1016/j.brainres.2011.11.043>
- 1525 Ellemers, N., Toorn, J. van der, Paunov, Y., & Leeuwen, T. van. (2019). The
1526 psychology of morality: A review and analysis of empirical studies published
1527 from 1940 through 2017. *Personality and Social Psychology Review*, 23(4),
1528 332–366. <https://doi.org/10.1177/1088868318811759>
- 1529 Enock, F. E., Hewstone, M. R. C., Lockwood, P. L., & Sui, J. (2020). Overlap in
1530 processing advantages for minimal ingroups and the self. *Scientific Reports*,
1531 10(1), 18933. <https://doi.org/10.1038/s41598-020-76001-9>
- 1532 Enock, F., Sui, J., Hewstone, M., & Humphreys, G. W. (2018). Self and team
1533 prioritisation effects in perceptual matching: Evidence for a shared
1534 representation. *Acta Psychologica*, 182, 107–118.
1535 <https://doi.org/10.1016/j.actpsy.2017.11.011>
- 1536 Faul, F., Erdfelder, E., Buchner, A., & Lang, A.-G. (2009). Statistical power
1537 analyses using g*power 3.1: Tests for correlation and regression analyses.
1538 *Behavior Research Methods*, 41(4), 1149–1160.
1539 <https://doi.org/10.3758/BRM.41.4.1149>
- 1540 Firestone, C., & Scholl, B. J. (2015). Enhanced visual awareness for morality and
1541 pajamas? Perception vs. Memory in ‘top-down’ effects. *Cognition*, 136, 409–416.
1542 <https://doi.org/10.1016/j.cognition.2014.10.014>
- 1543 Freeman, J. B., & Ambady, N. (2011). A dynamic interactive theory of person
1544 construal. *Psychological Review*, 118(2), 247–279.

- 1545 https://doi.org/10.1037/a0022327
- 1546 Gantman, A. P., & Van Bavel, J. J. (2014). The moral pop-out effect: Enhanced
1547 perceptual awareness of morally relevant stimuli. *Cognition*, 132(1), 22–29.
1548 https://doi.org/10.1016/j.cognition.2014.02.007
- 1549 Goh, J. X., Hall, J. A., & Rosenthal, R. (2016). Mini meta-analysis of your own
1550 studies: Some arguments on why and a primer on how [Journal Article]. *Social*
1551 and *Personality Psychology Compass*, 10(10), 535–549.
1552 https://doi.org/10.1111/spc.12267
- 1553 Goodwin, G. P. (2015). Moral character in person perception. *Current Directions in*
1554 *Psychological Science*, 24(1), 38–44. https://doi.org/10.1177/0963721414550709
- 1555 Goodwin, G. P., Piazza, J., & Rozin, P. (2014). Moral character predominates in
1556 person perception and evaluation. *Journal of Personality and Social Psychology*,
1557 106(1), 148–168. https://doi.org/10.1037/a0034726
- 1558 Henrich, J., Heine, S. J., & Norenzayan, A. (2010). The weirdest people in the
1559 world? *Behavioral and Brain Sciences*, 33(2), 61–83.
1560 https://doi.org/10.1017/S0140525X0999152X
- 1561 Hofmann, W., Wisneski, D. C., Brandt, M. J., & Skitka, L. J. (2014). Morality in
1562 everyday life. *Science*, 345(6202), 1340–1343.
1563 https://doi.org/10.1126/science.1251560
- 1564 Hu, C.-P., Lan, Y., Macrae, C. N., & Sui, J. (2020). Good me bad me: Does valence
1565 influence self-prioritization during perceptual decision-making? [Journal Article].
1566 *Collabra: Psychology*, 6(1), 20. https://doi.org/10.1525/collabra.301
- 1567 Jones, B. C., DeBruine, L. M., Flake, J. K., Liuzza, M. T., Antfolk, J., Arinze, N.
1568 C., ... Coles, N. A. (2020). To which world regions does the valence-dominance
1569 model of social perception apply? *Nature Human Behaviour*.

- 1570 <https://doi.org/10.31234/osf.io/n26dy>
- 1571 Juechems, K., & Summerfield, C. (2019). Where does value come from? *Trends in*
1572 *Cognitive Sciences*, 23(10), 836–850. <https://doi.org/10.1016/j.tics.2019.07.012>
- 1573 Kirby, K. N., & Gerlanc, D. (2013). BootES: An r package for bootstrap confidence
1574 intervals on effect sizes. *Behavior Research Methods*, 45(4), 905–927.
1575 <https://doi.org/10.3758/s13428-013-0330-5>
- 1576 Klein, N., & Epley, N. (2016). Maybe holier, but definitely less evil, than you:
1577 Bounded self-righteousness in social judgment. *Journal of Personality and Social*
1578 *Psychology*, 110(5), 660–674. <https://doi.org/10.1037/pspa0000050>
- 1579 Liu, Q., Wang, F., Yan, W., Peng, K., Sui, J., & Hu, C.-P. (2020). Questionnaire
1580 data from the revision of a chinese version of free will and determinism plus
1581 scale [Journal Article]. *Journal of Open Psychology Data*, 8(1), 1.
1582 <https://doi.org/10.5334/jopd.49/>
- 1583 Matzke, D., & Wagenmakers, E.-J. (2009). Psychological interpretation of the
1584 ex-gaussian and shifted wald parameters: A diffusion model analysis.
1585 *Psychonomic Bulletin & Review*, 16(5), 798–817.
1586 <https://doi.org/10.3758/PBR.16.5.798>
- 1587 McHugh, C., McGann, M., Igou, E. R., & Kinsella, E. (2019). *Moral judgment as*
1588 *categorization (MJAC)*. PsyArXiv. <https://doi.org/10.31234/osf.io/72dzp>
- 1589 McNeish, D., & Wolf, M. G. (2020). Thinking twice about sum scores. *Behavior*
1590 *Research Methods*. <https://doi.org/10.3758/s13428-020-01398-0>
- 1591 Monin, B., & Jordan, A. H. (2009). The dynamic moral self: A social psychological
1592 perspective. In *Personality, identity, and character: Explorations in moral*
1593 *psychology* (pp. 341–354). New York, NY, US: Cambridge University Press.
1594 <https://doi.org/10.1017/CBO9780511627125.016>

- 1595 Pelli, D. G. (1997). The VideoToolbox software for visual psychophysics:
1596 Transforming numbers into movies [Journal Article]. *Spatial Vision*, 10(4),
1597 437–442.
- 1598 Reicher, S., & Hopkins, N. (2016). Perception, action, and the social dynamics of
1599 the variable self. *Psychological Inquiry*, 27(4), 341–347.
1600 <https://doi.org/10.1080/1047840X.2016.1217584>
- 1601 Rouder, J. N., & Lu, J. (2005). An introduction to bayesian hierarchical models
1602 with an application in the theory of signal detection [Journal Article].
1603 *Psychonomic Bulletin & Review*, 12(4), 573–604.
1604 <https://doi.org/10.3758/bf03196750>
- 1605 Rousselet, G. A., & Wilcox, R. R. (2019). Reaction times and other skewed
1606 distributions: Problems with the mean and the median [Preprint].
1607 *Meta-Psychology*. <https://doi.org/10.1101/383935>
- 1608 Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2013). *Life after p-hacking*
1609 [Conference Proceedings]. <https://doi.org/10.2139/ssrn.2205186>
- 1610 Stanley, M. L., & De Brigard, F. (2019). Moral memories and the belief in the good
1611 self. *Current Directions in Psychological Science*, 28(4), 387–391.
1612 <https://doi.org/10.1177/0963721419847990>
- 1613 Stein, T., Grubb, C., Bertrand, M., Suh, S. M., & Verosky, S. C. (2017). No impact
1614 of affective person knowledge on visual awareness: Evidence from binocular
1615 rivalry and continuous flash suppression. *Emotion*, 17(8), 1199–1207.
1616 <https://doi.org/10.1037/emo0000305>
- 1617 Stolier, R. M., & Freeman, J. B. (2016). Functional and temporal considerations for
1618 top-down influences in social perception. *Psychological Inquiry*, 27(4), 352–357.
1619 <https://doi.org/10.1080/1047840X.2016.1216034>

- 1620 Strohminger, N., Knobe, J., & Newman, G. (2017). The true self: A psychological
1621 concept distinct from the self: *Perspectives on Psychological Science*.
1622 <https://doi.org/10.1177/1745691616689495>
- 1623 Sui, J., He, X., & Humphreys, G. W. (2012). Perceptual effects of social salience:
1624 Evidence from self-prioritization effects on perceptual matching [Journal
1625 Article]. *Journal of Experimental Psychology: Human Perception and*
1626 *Performance*, 38(5), 1105–1117. <https://doi.org/10.1037/a0029792>
- 1627 Tappin, B. M., & McKay, R. T. (2017). The illusion of moral superiority. *Social*
1628 *Psychological and Personality Science*, 8(6), 623–631.
1629 <https://doi.org/10.1177/1948550616673878>
- 1630 Turner, J. C., Hogg, M. A., Oakes, P. J., Reicher, S. D., & Wetherell, M. S. (1987).
1631 *Rediscovering the social group: A self-categorization theory*. Cambridge, MA,
1632 US: Basil Blackwell.
- 1633 Turner, J. C., Oakes, P. J., Haslam, S. A., & McGarty, C. (1994). Self and
1634 collective: Cognition and social context. *Personality and Social Psychology*
1635 *Bulletin*, 20(5), 454–463. <https://doi.org/10.1177/0146167294205002>
- 1636 Uhlmann, E. L., Pizarro, D. A., & Diermeier, D. (2015). A person-centered
1637 approach to moral judgment: *Perspectives on Psychological Science*.
1638 <https://doi.org/10.1177/1745691614556679>
- 1639 Verosky, S. C., & Todorov, A. (2010). Differential neural responses to faces
1640 physically similar to the self as a function of their valence. *NeuroImage*, 49(2),
1641 1690–1698. <https://doi.org/10.1016/j.neuroimage.2009.10.017>
- 1642 Waytz, A., Dungan, J., & Young, L. (2013). The whistleblower's dilemma and the
1643 fairness–loyalty tradeoff. *Journal of Experimental Social Psychology*, 49(6),
1644 1027–1033. <https://doi.org/10.1016/j.jesp.2013.07.002>

- 1645 Wiecki, T. V., Sofer, I., & Frank, M. J. (2013). HDDM: Hierarchical bayesian
1646 estimation of the drift-diffusion model in python. *Frontiers in Neuroinformatics*,
1647 7. <https://doi.org/10.3389/fninf.2013.00014>
- 1648 Willis, J., & Todorov, A. (2006). First impressions: Making up your mind after a
1649 100-ms exposure to a face. *Psychological Science*, 17(7), 592–598.
1650 <https://doi.org/10.1111/j.1467-9280.2006.01750.x>
- 1651 Xiao, Y. J., Coppin, G., & Bavel, J. J. V. (2016). Perceiving the world through
1652 group-colored glasses: A perceptual model of intergroup relations. *Psychological
1653 Inquiry*, 27(4), 255–274. <https://doi.org/10.1080/1047840X.2016.1199221>

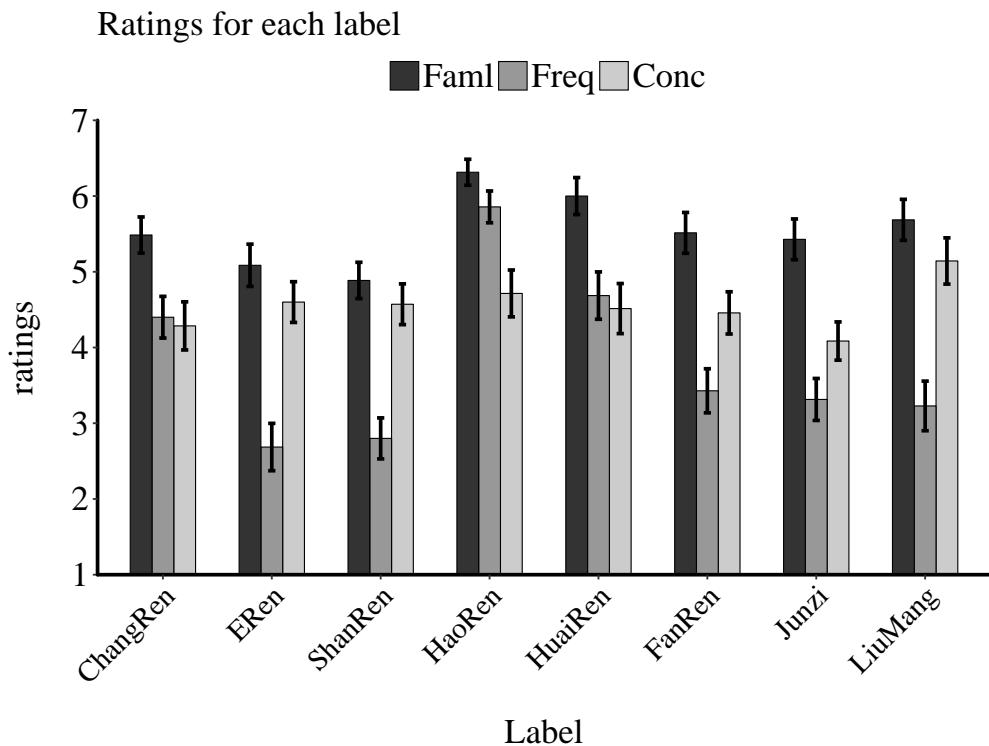


Figure 1. Ratings of words in exp 1b

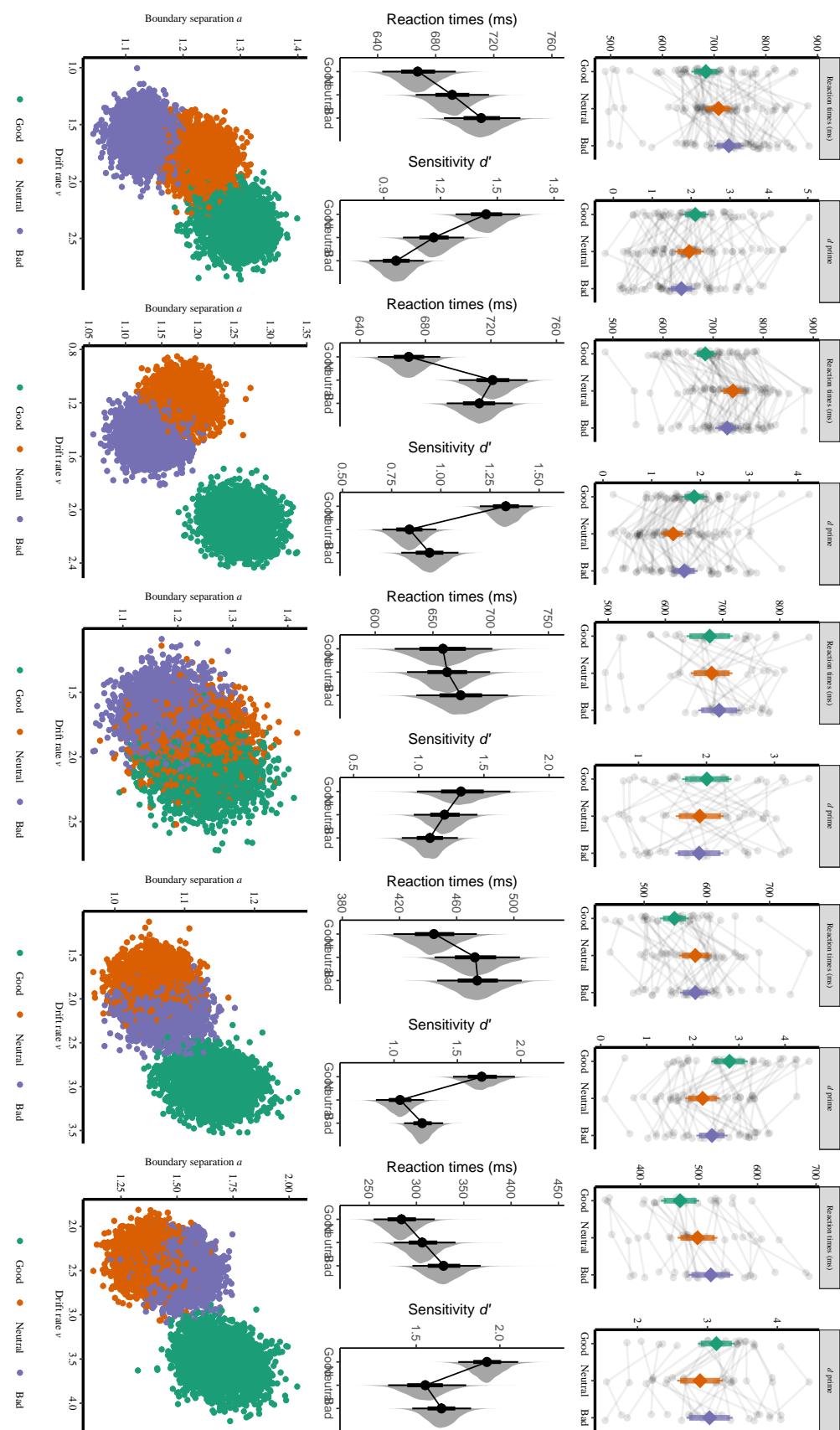


Figure 2. Results for part 1.

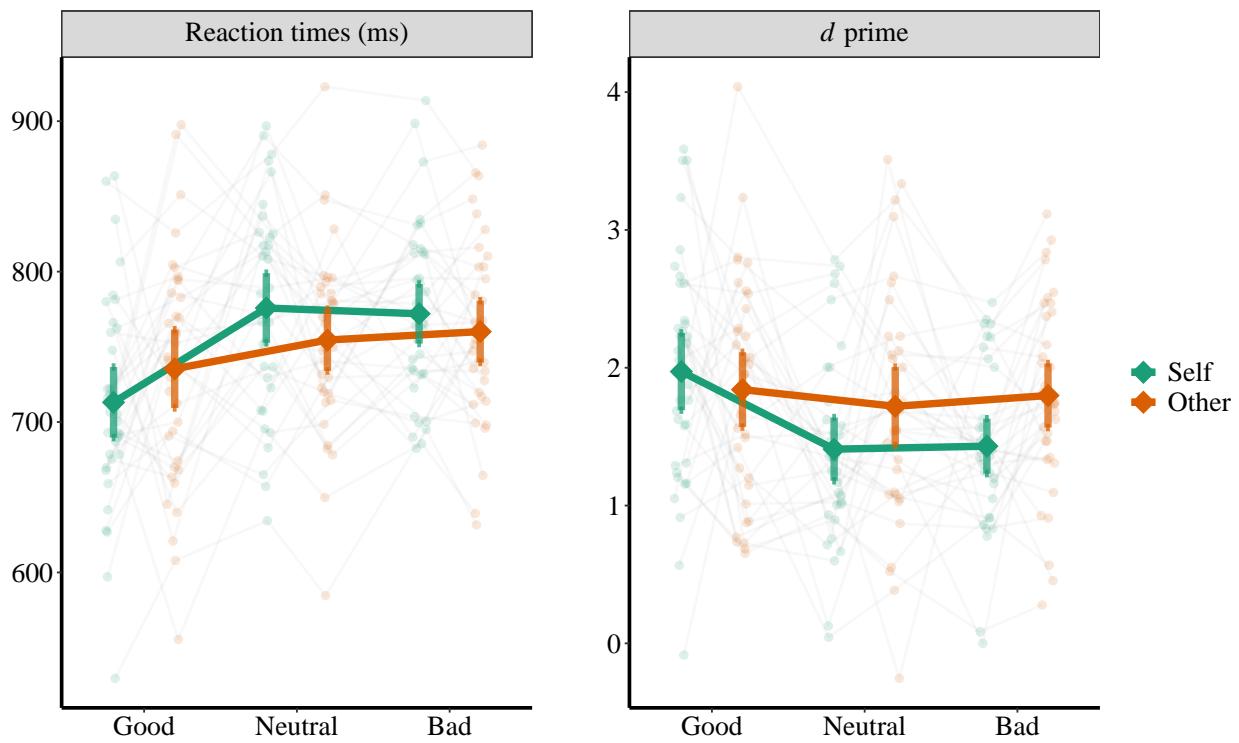


Figure 3. RT and d' of Experiment 3a.

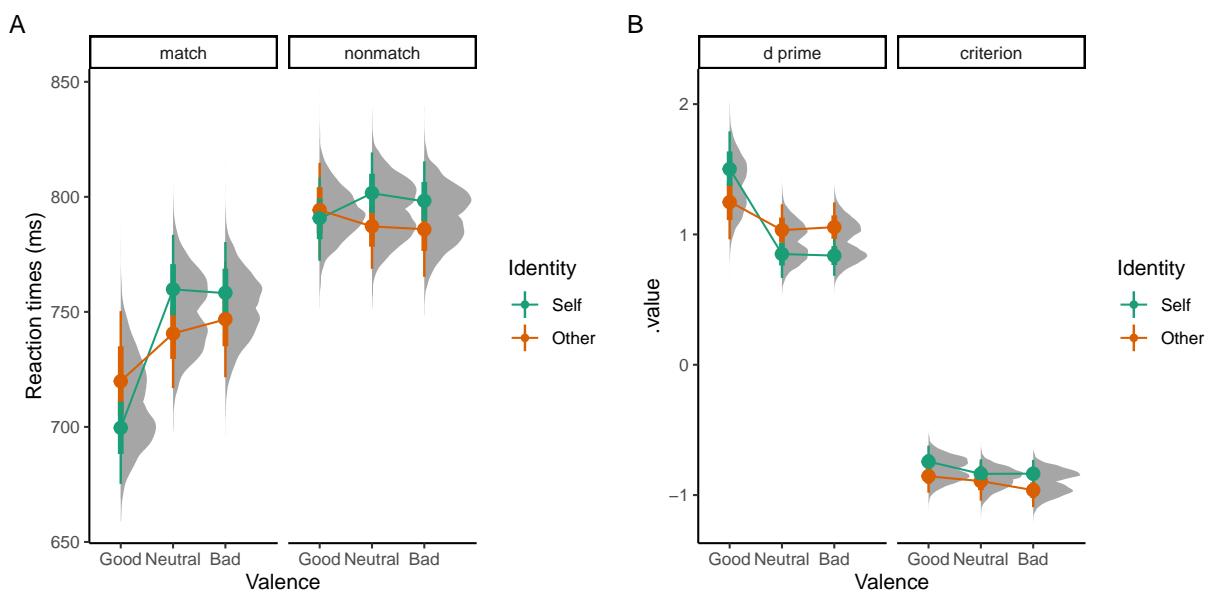


Figure 4. Exp3a: Results of Bayesian GLM analysis.

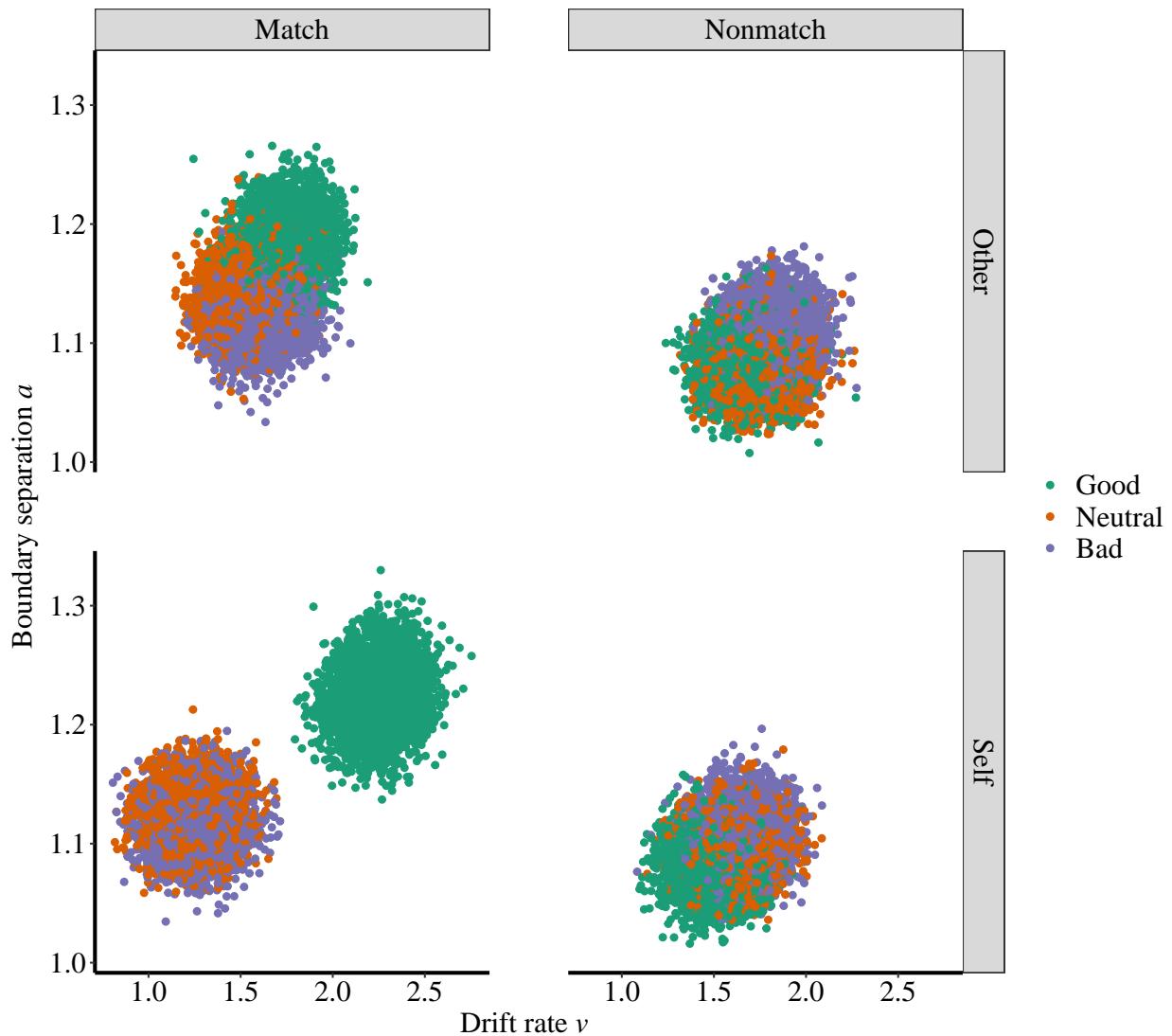


Figure 5. Exp3a: Results of HDDM.

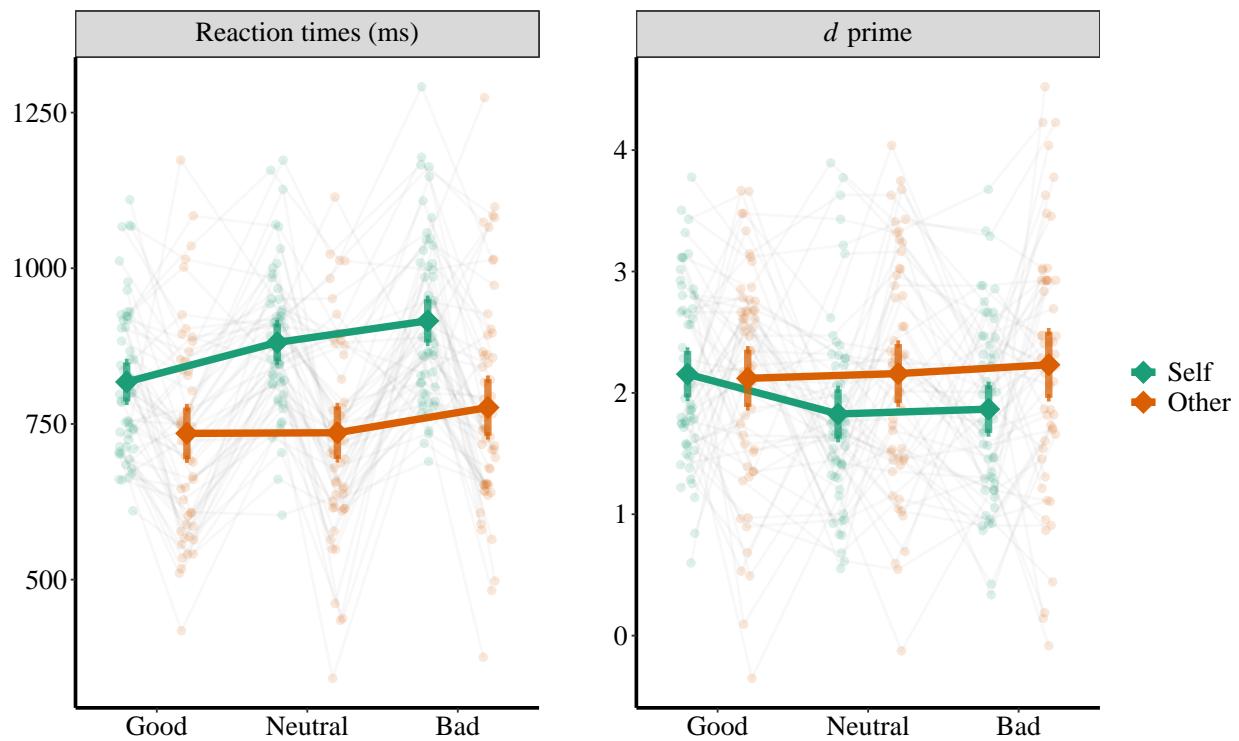


Figure 6. RT and d' of Experiment 3b.

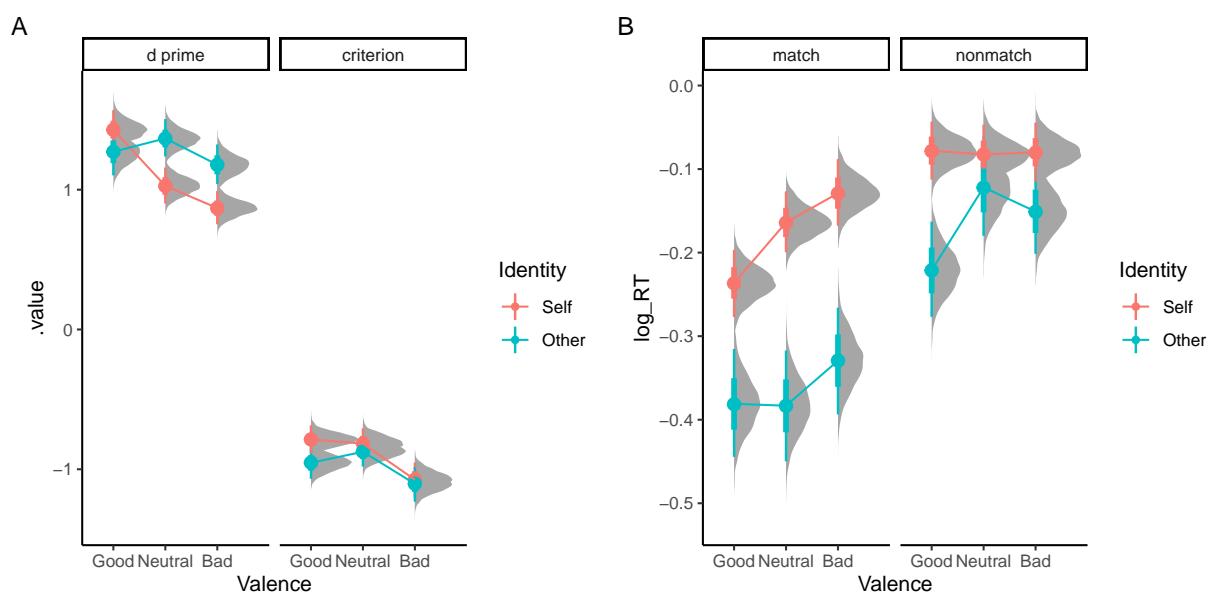


Figure 7. exp3b: Results of Bayesian GLM analysis.

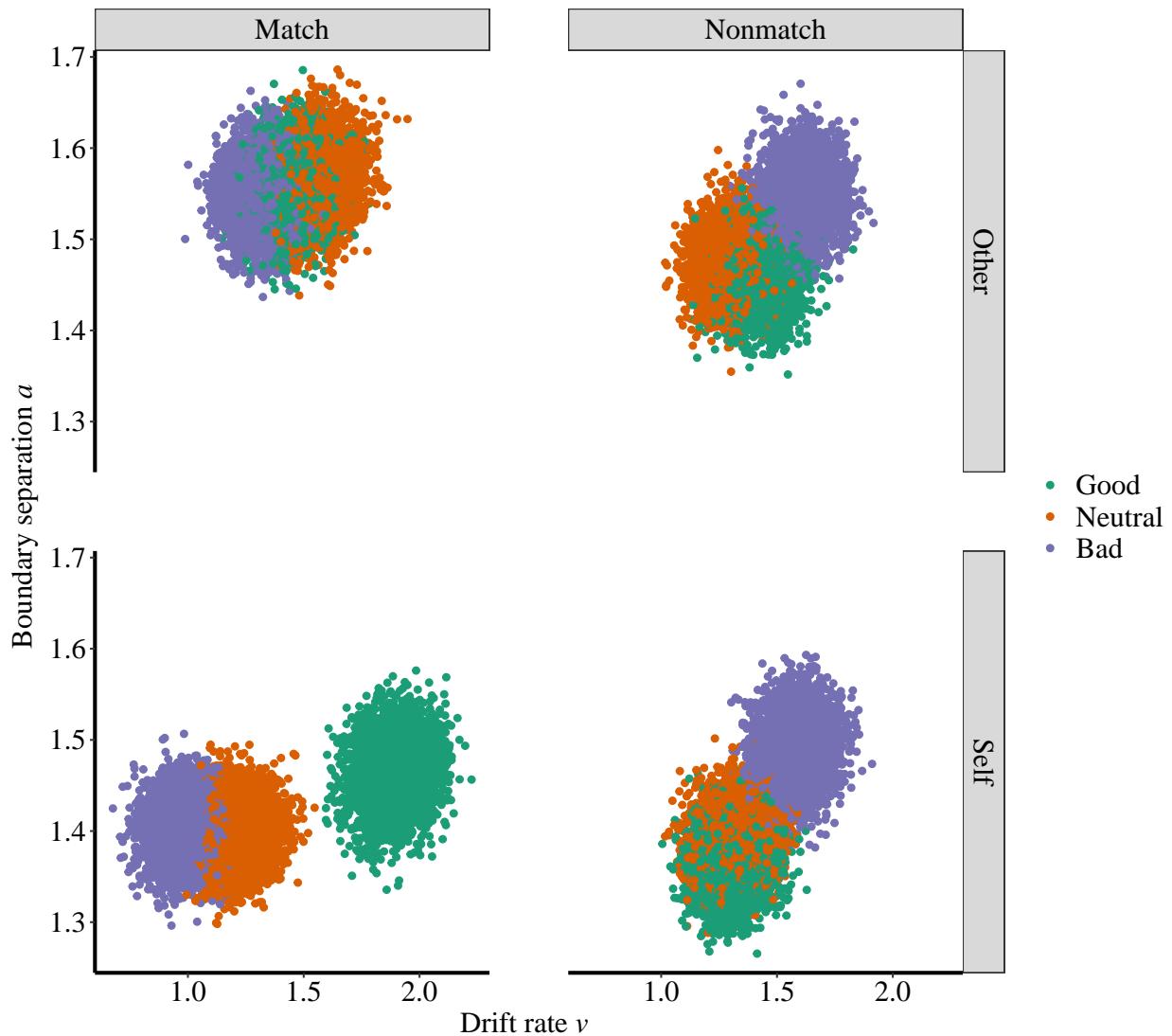


Figure 8. exp3b: Results of HDDM.

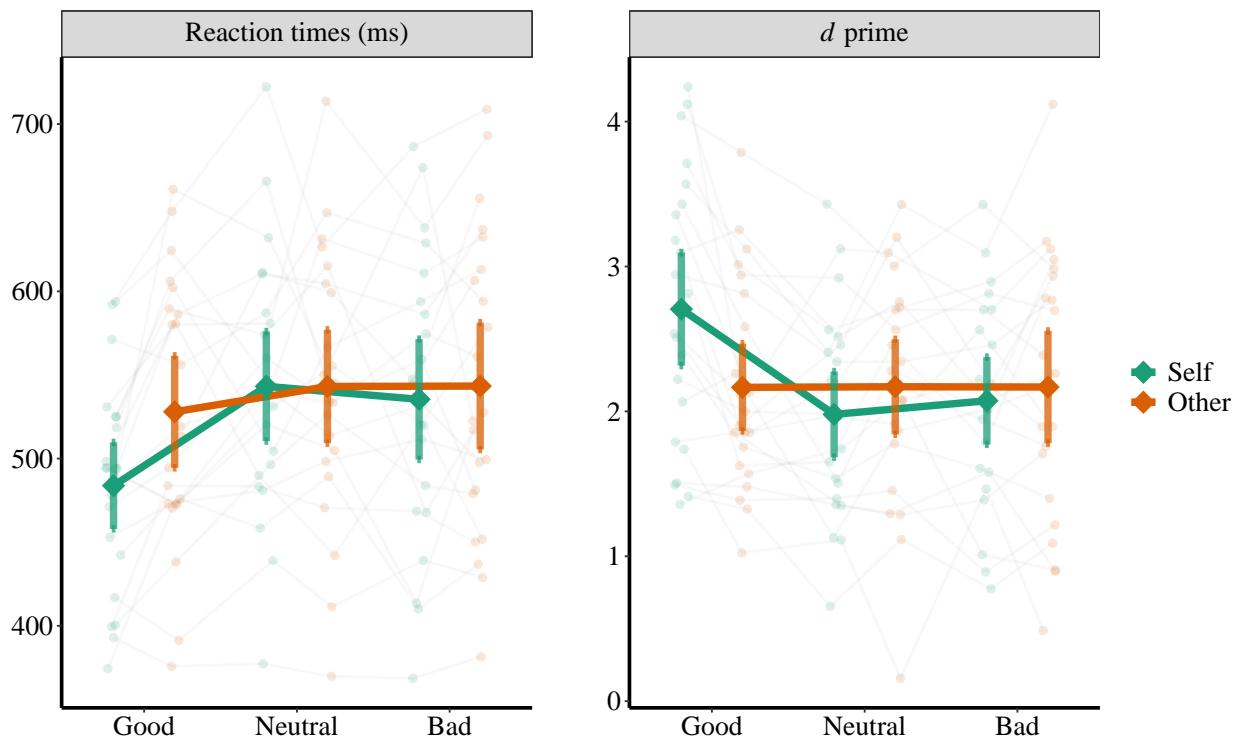


Figure 9. RT and d' of Experiment 6b.

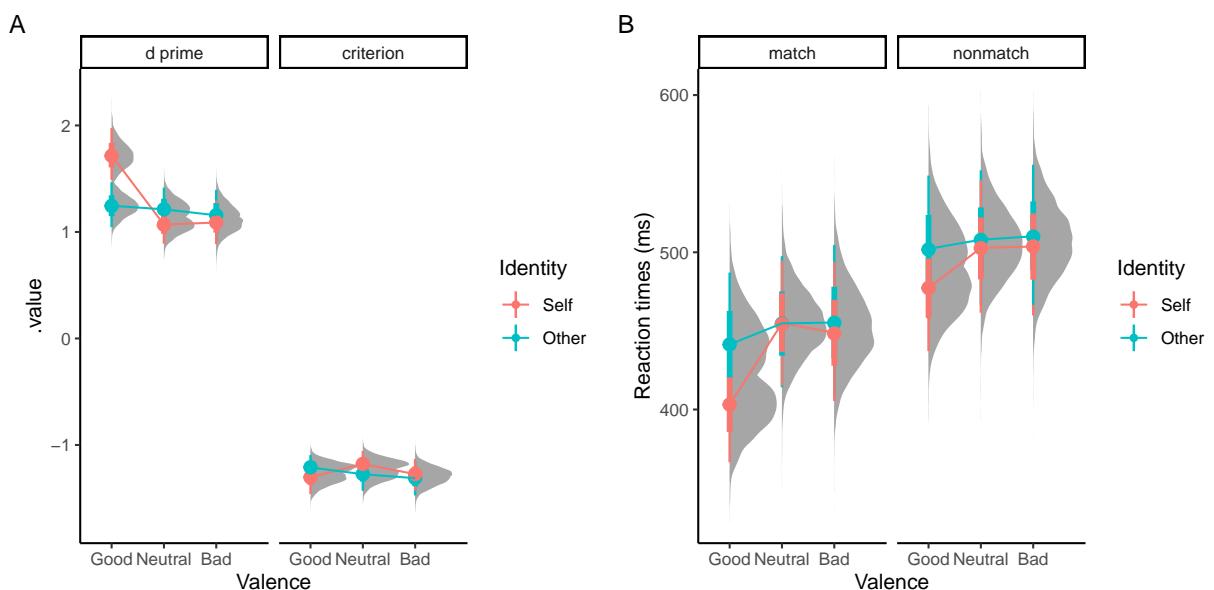


Figure 10. exp6b_d1: Results of Bayesian GLM analysis.

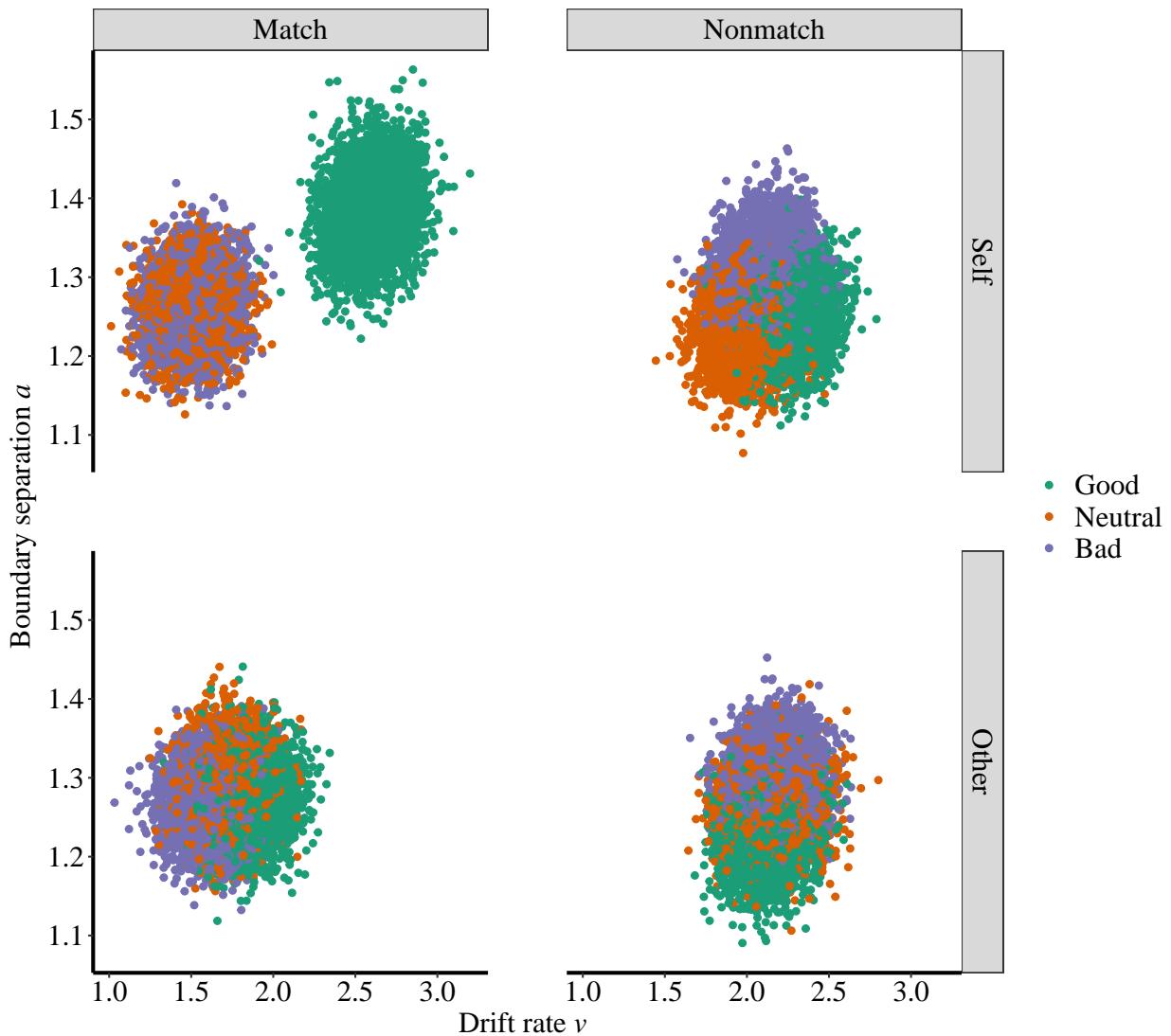


Figure 11. exp6b: Results of HDDM (Day 1).

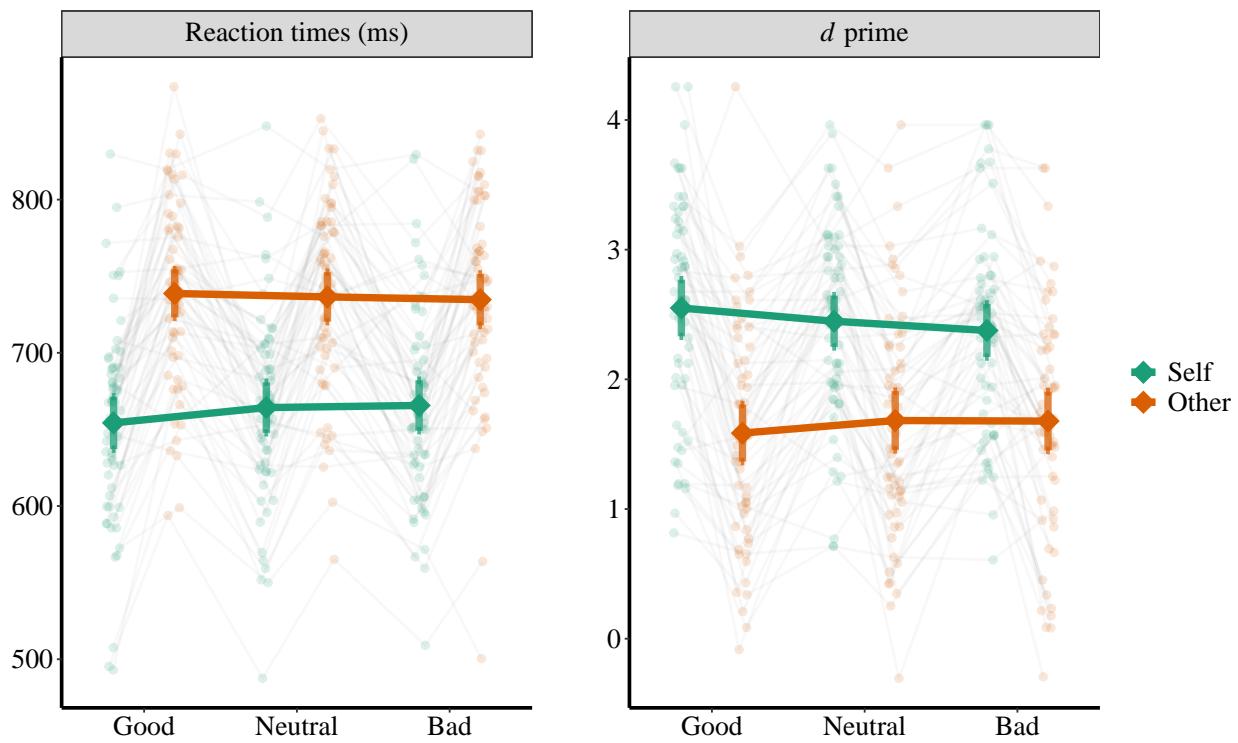


Figure 12. RT and d' of Experiment 4a.

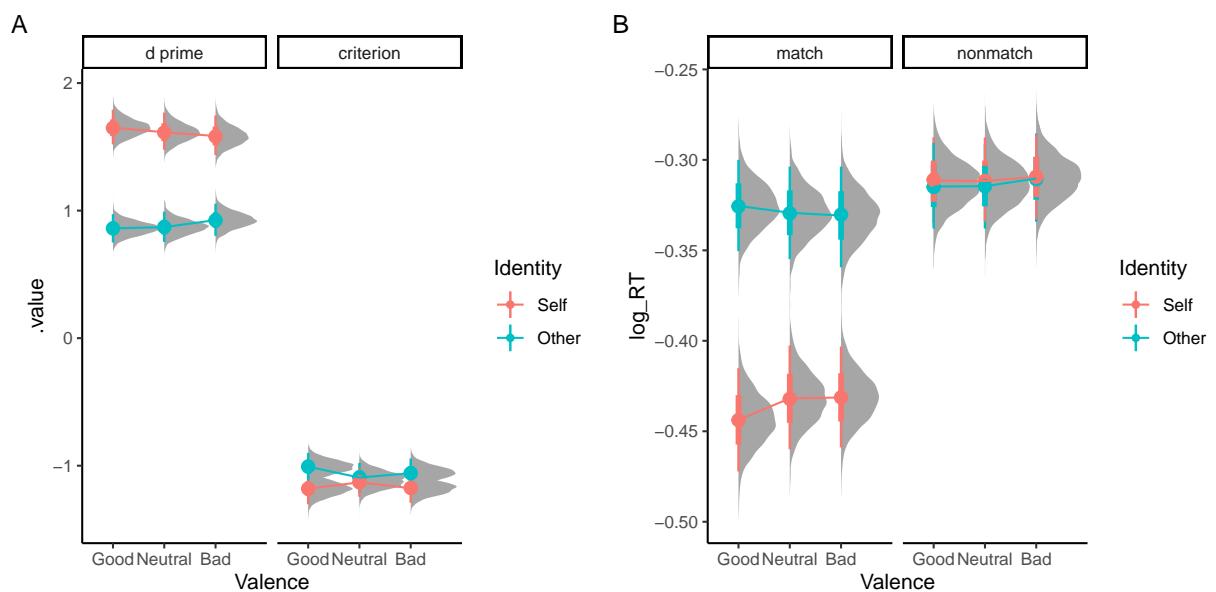


Figure 13. exp4a: Results of Bayesian GLM analysis.

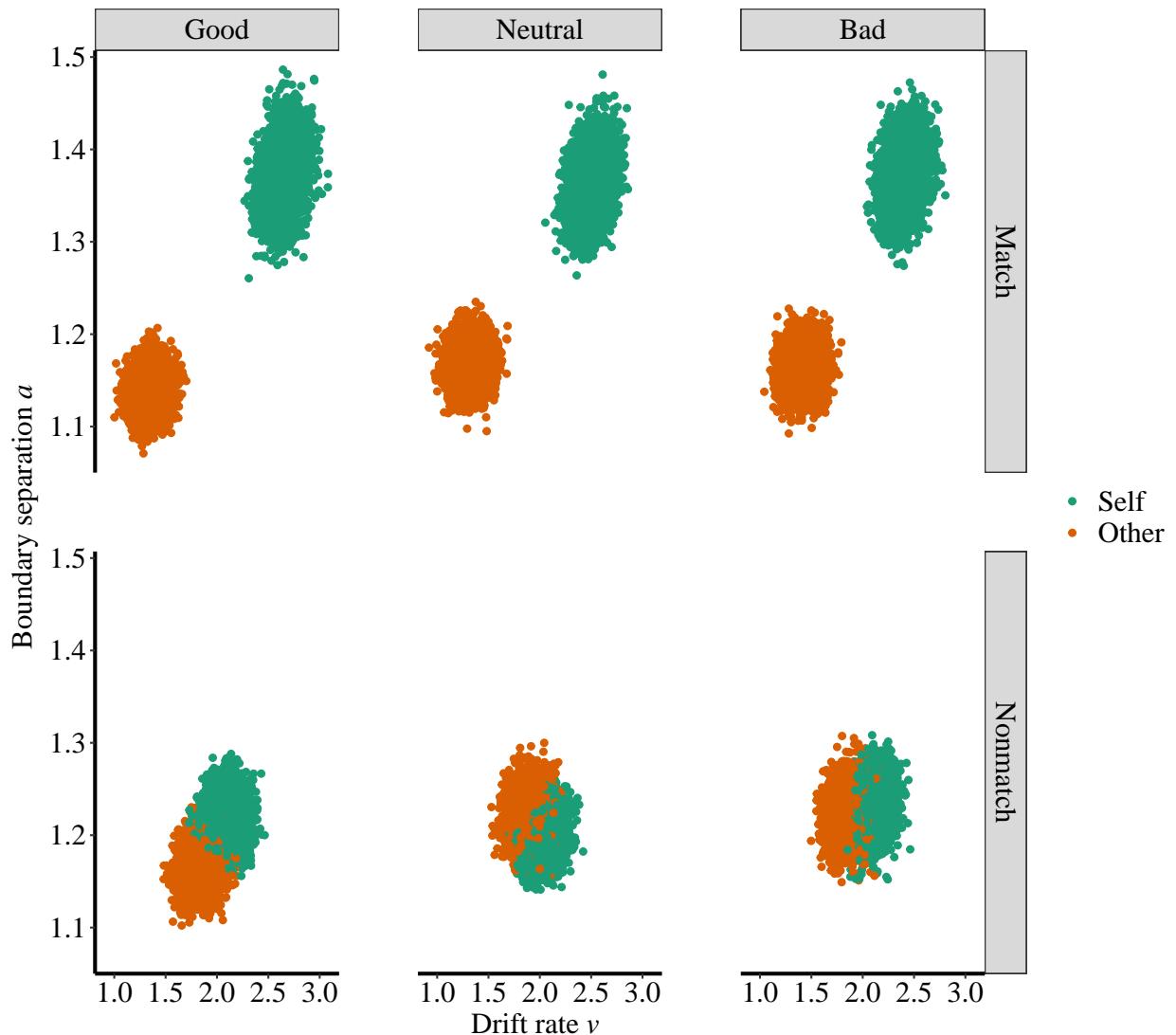


Figure 14. exp4a: Results of HDDM.

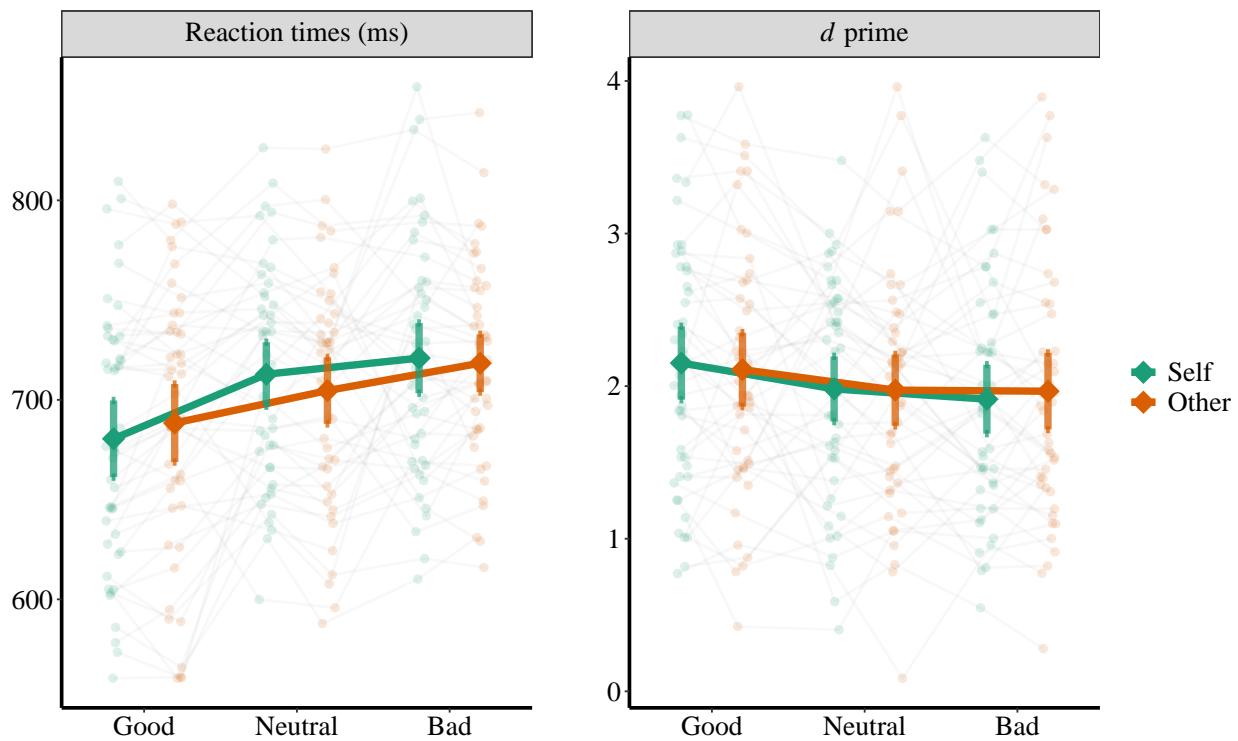


Figure 15. RT and d' prime of Experiment 4b.

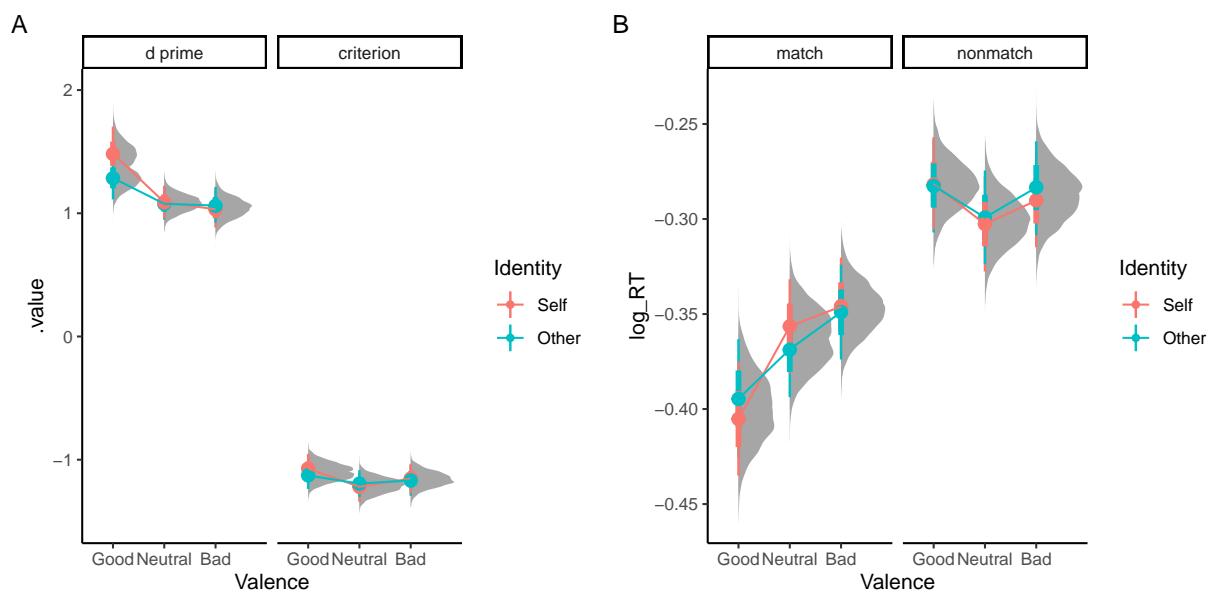


Figure 16. exp4b: Results of Bayesian GLM analysis.

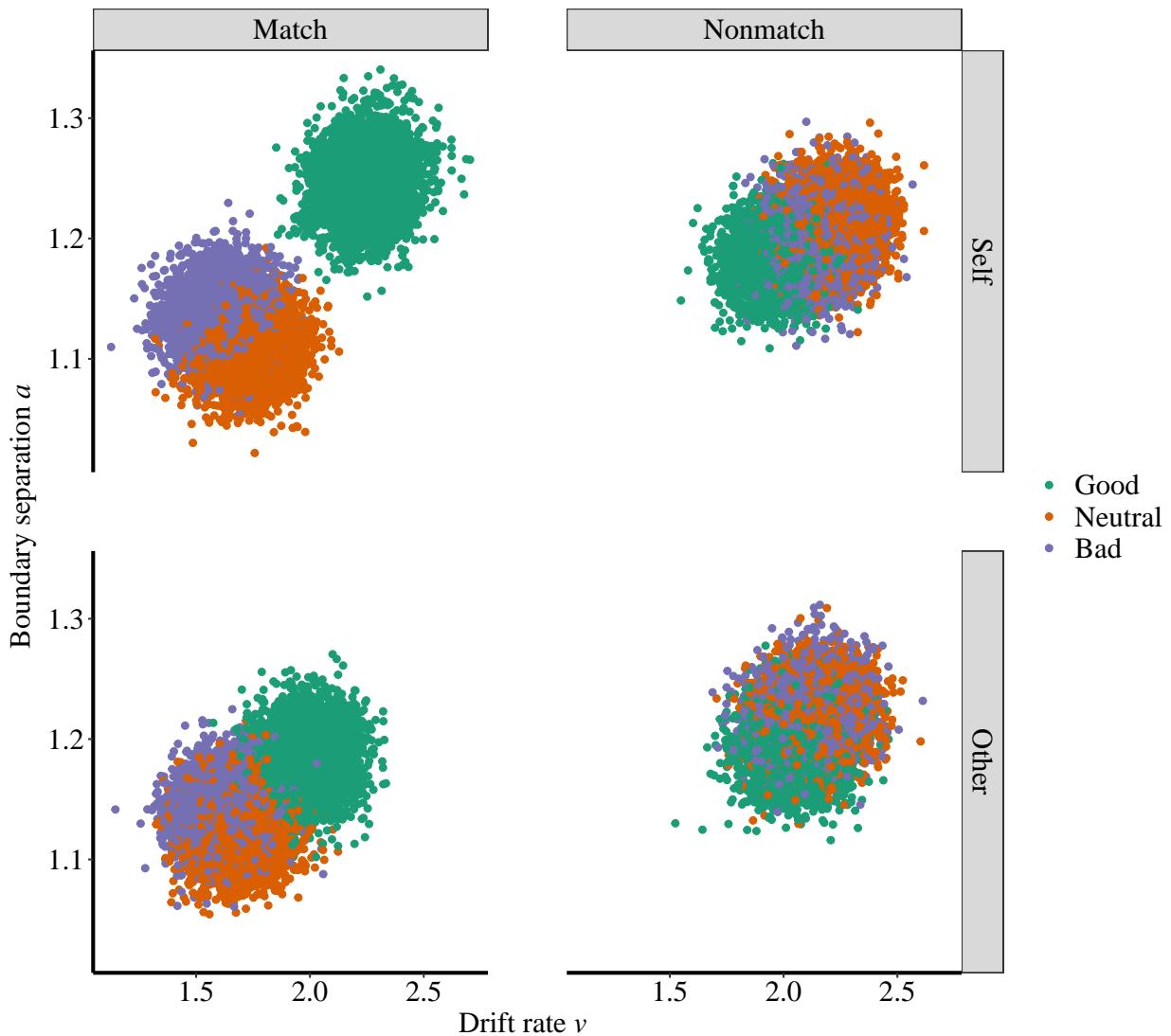


Figure 17. exp4b: Results of HDDM.

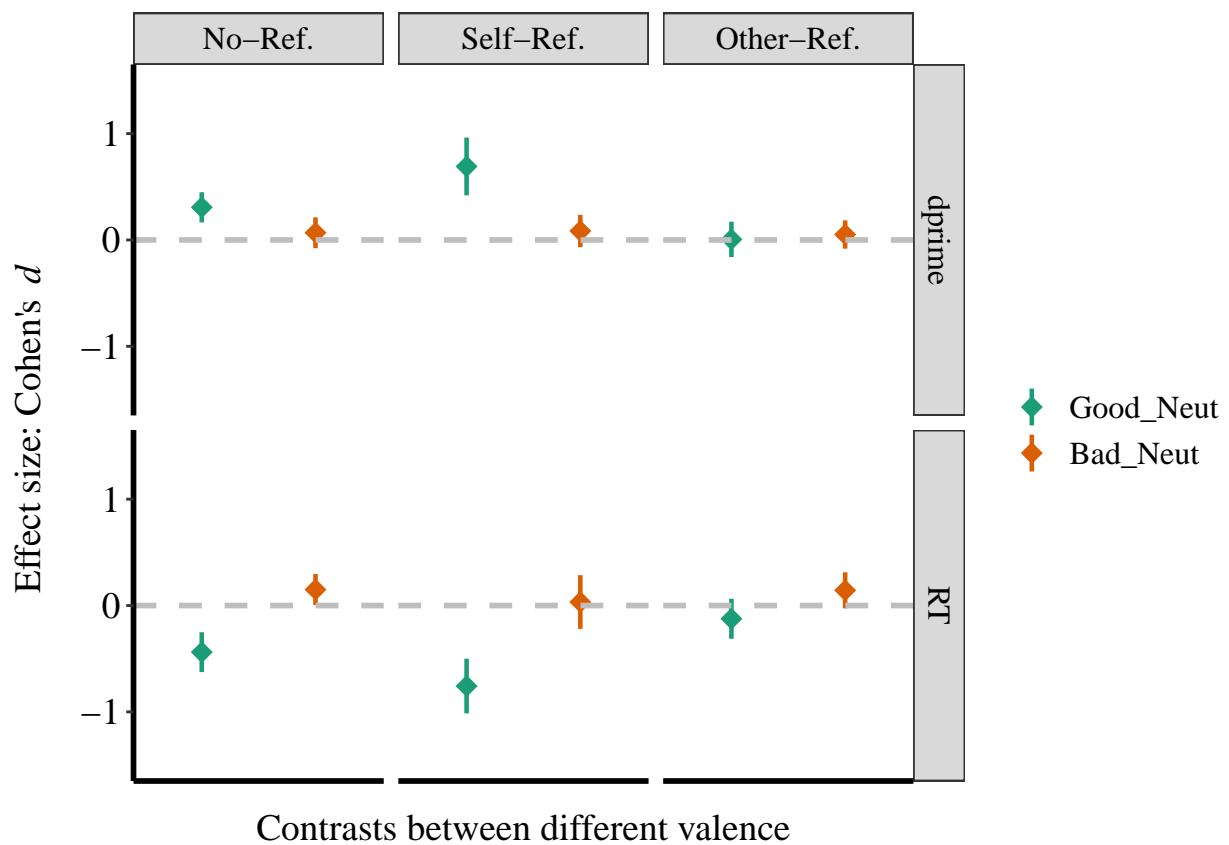


Figure 18. Effect size (Cohen's d) of Valence.

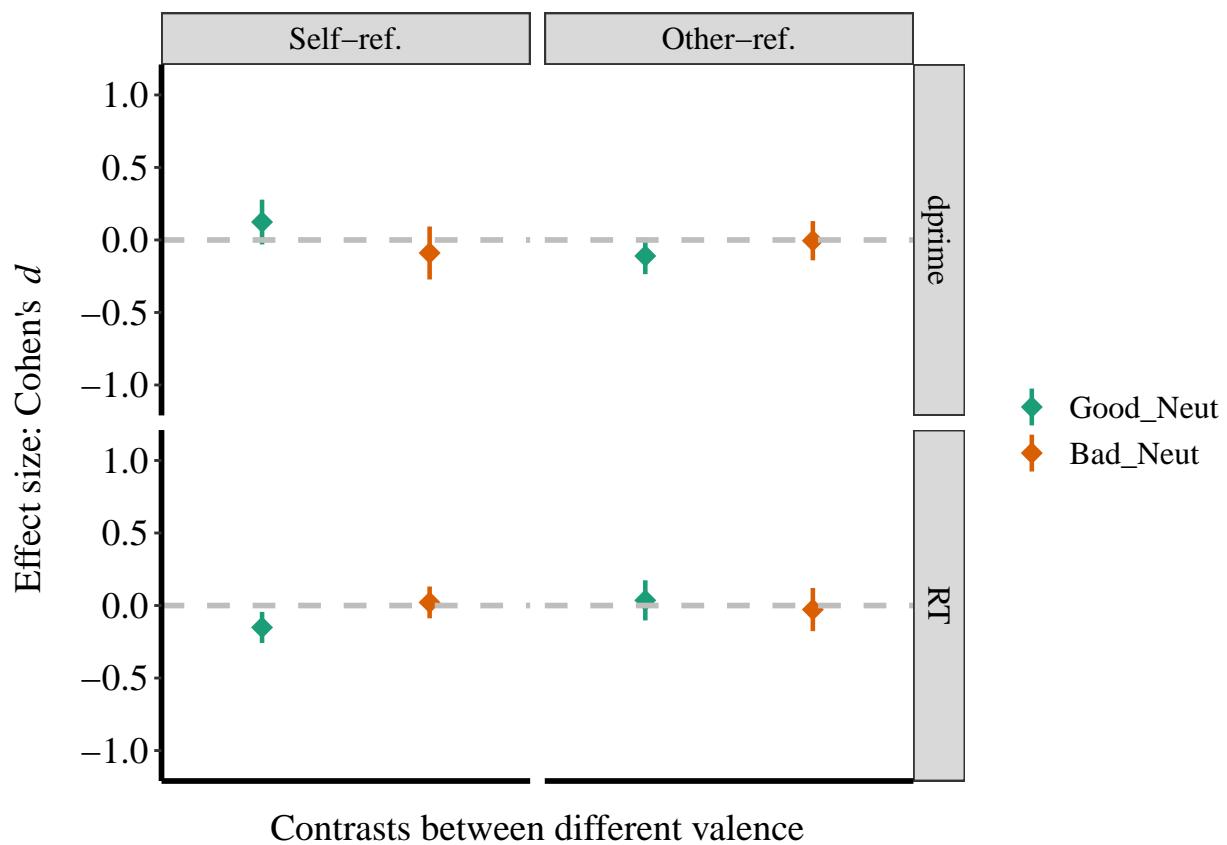


Figure 19. Effect size (Cohen's d) of Valence in Exp4a.

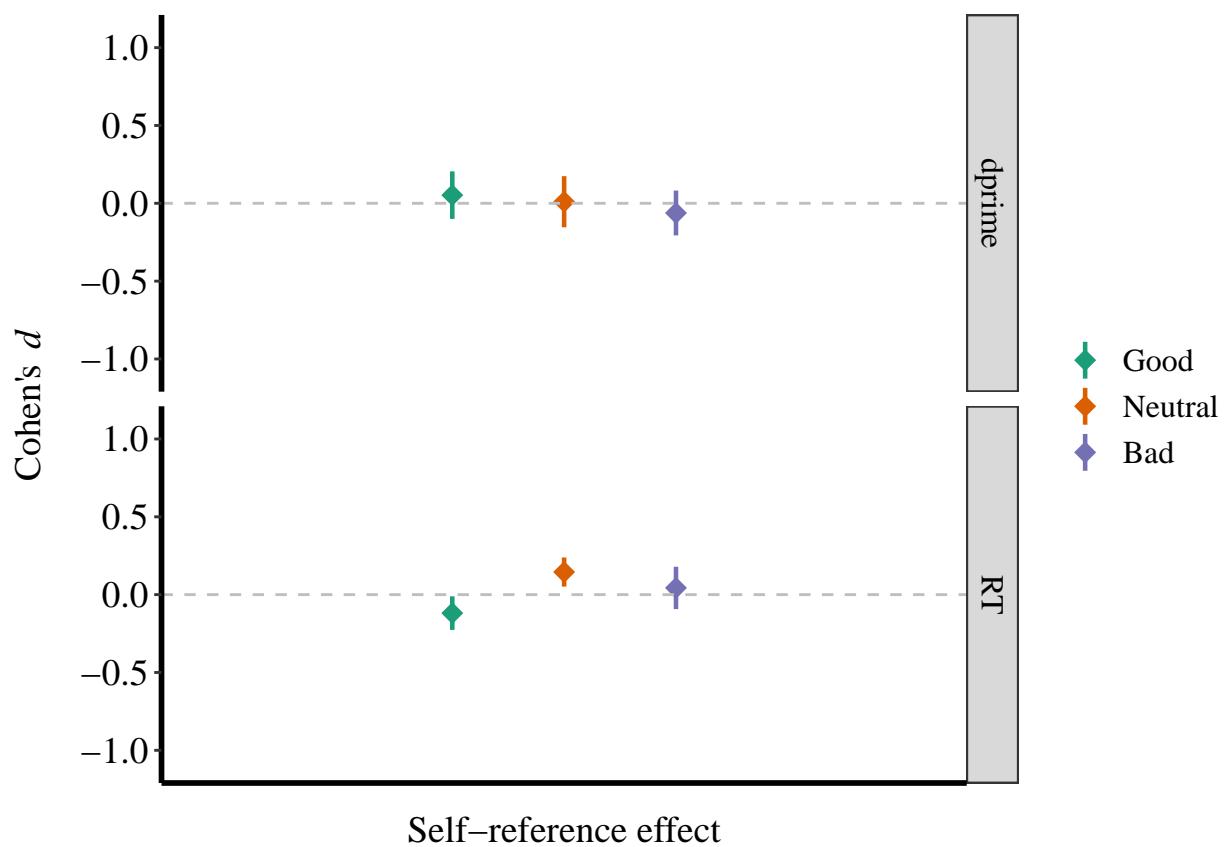


Figure 20. Effect size (Cohen's d) of Valence in Exp4b.

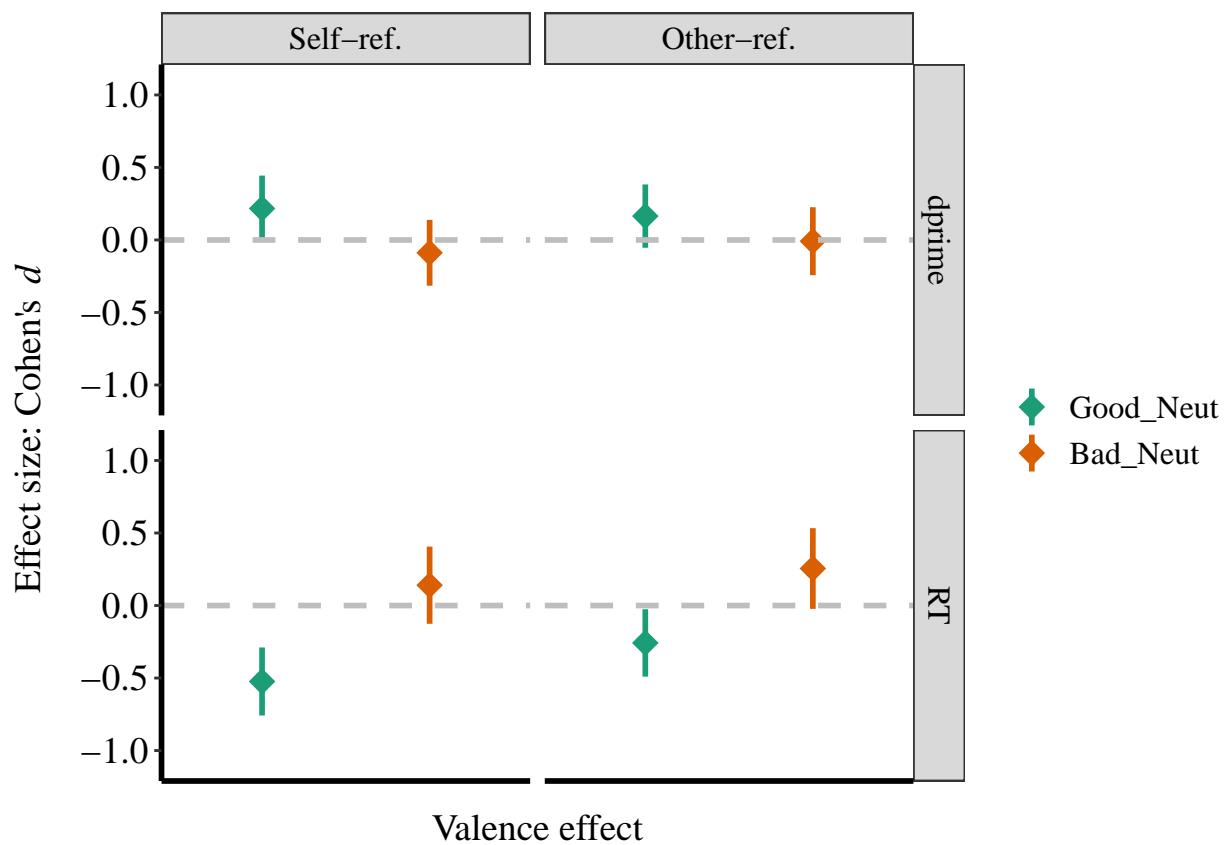


Figure 21. Effect size (Cohen's d) of Valence in Exp4b.

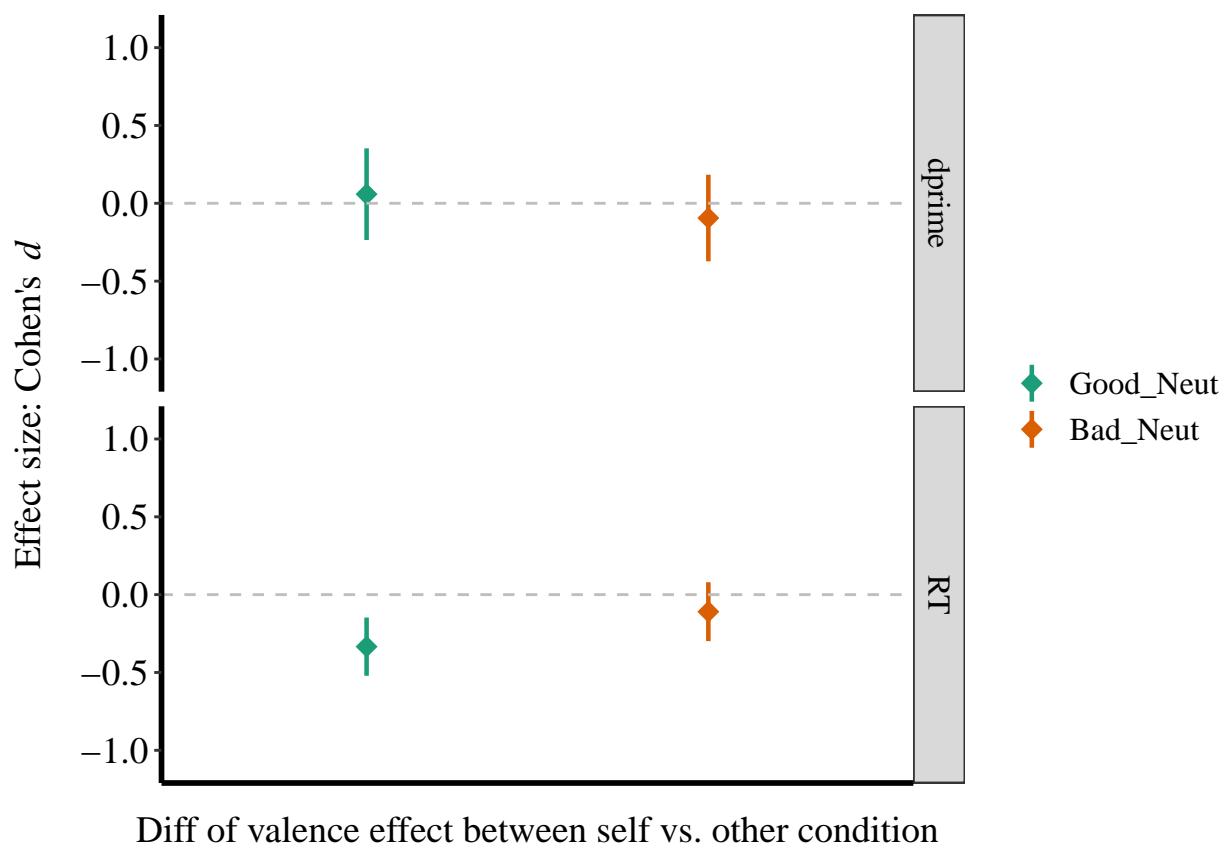


Figure 22. Effect size (Cohen's d) of Valence in Exp4b.

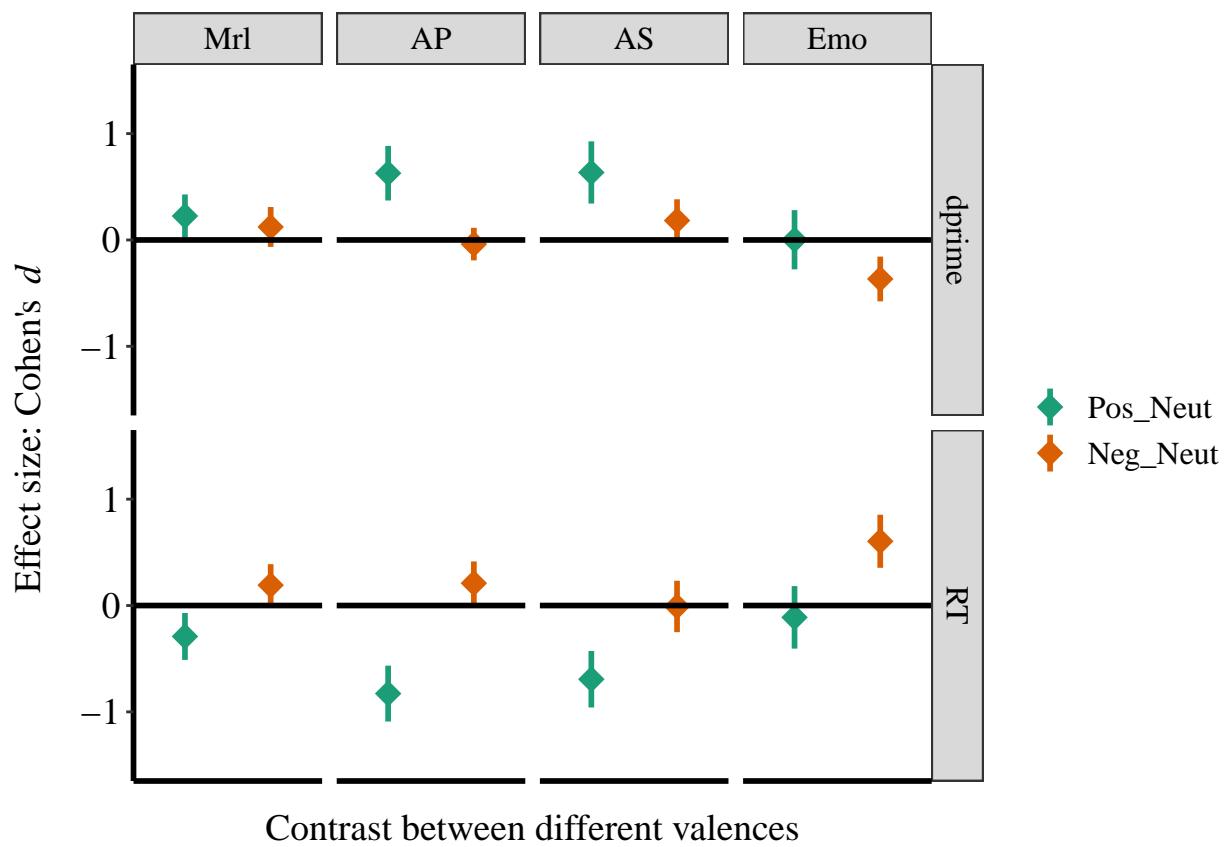


Figure 23. Effect size (Cohen's d) of Valence in Exp5.

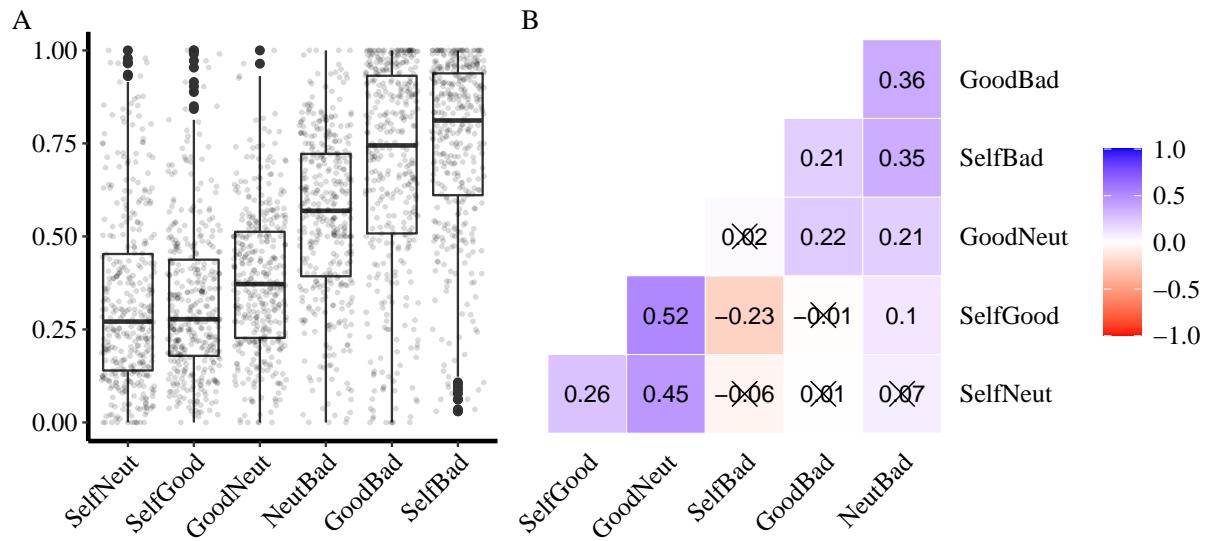


Figure 24. Self-rated personal distance

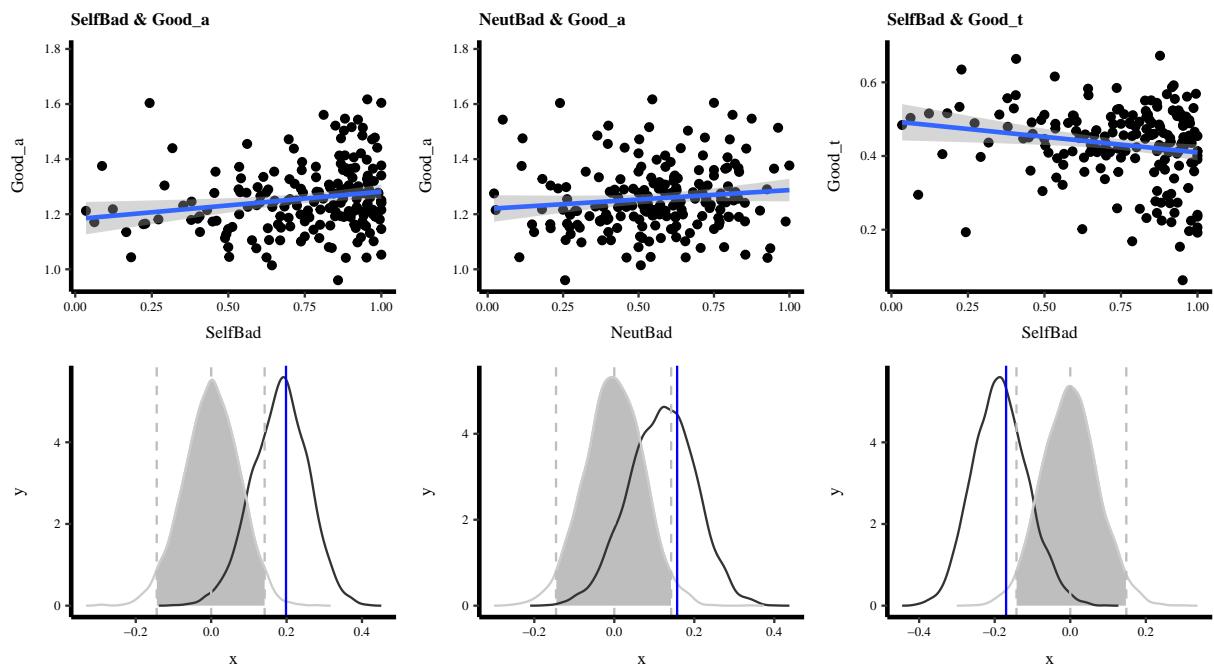


Figure 25. Correlation between moral identity and boundary separation of good condition; moral self-image and drift rate of good condition

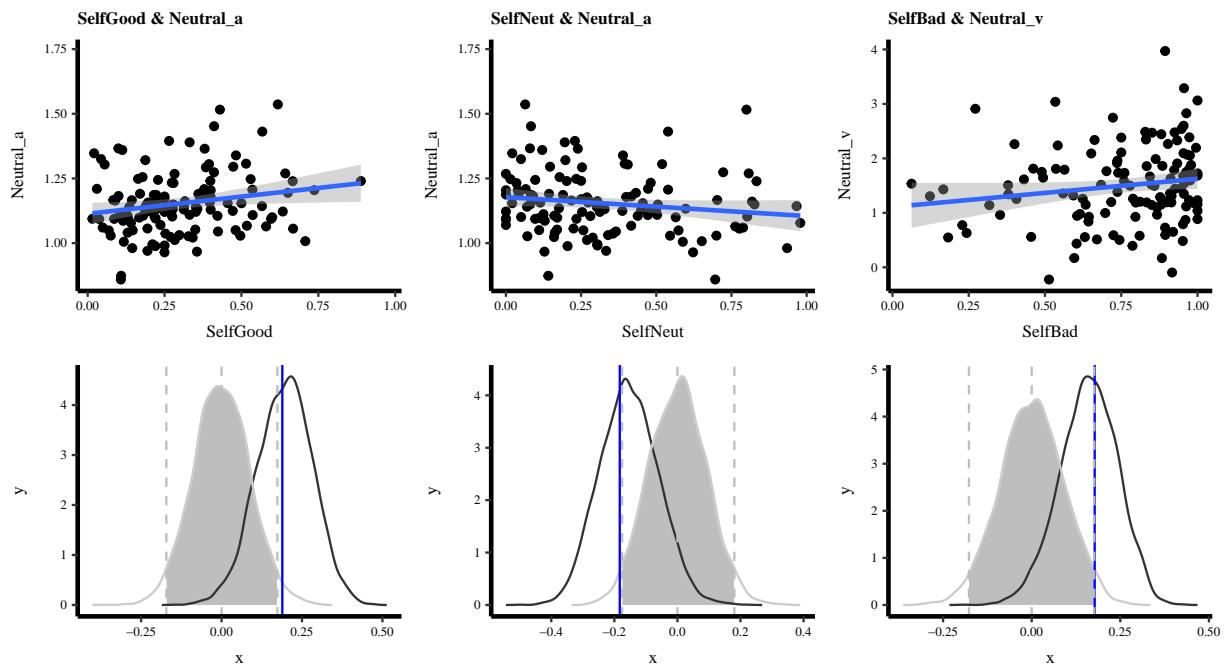


Figure 26. Correlation between personal distance and boundary separation of neutral condition