

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

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

³ ¹ TBA

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

⁵ ³ Tsinghua University, 100084 Beijing, China

⁶ ⁴ University of Aberdeen, Aberdeen, Scotland

⁷ Author Note

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

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

¹⁰ Psychology, University of Aberdeen, Aberdeen, Scotland.

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

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

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

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

¹⁵ Germany. E-mail: hcp4715@gmail.com

16

Abstract

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

37

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

38

Word count: X

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

40 **Introduction**

41 [sentences in bracket are key ideas]

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

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

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

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

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

83 Perceptual process of moral character

84 [exp1a, b, c, and exp2]

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

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

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

¹⁰⁴ An alternative is to use abstract semantic concepts. Abstract concepts of moral
¹⁰⁵ character are used to describe and represent moral characters. These abstract concepts
¹⁰⁶ may be part of a dynamic network in which sensory cue, concrete behaviors and other
¹⁰⁷ information can activate/inhibit each other (e.g., aggressiveness) (Amodio, 2019; Freeman
¹⁰⁸ & Ambady, 2011). If a concept of moral character (e.g., good person) is activated, it
¹⁰⁹ should be able to influence on the perceptual process of the visual cues through the
¹¹⁰ dynamic network, especially when the perceptual decision-making is about the concept-cue
¹¹¹ association. In this case, abstract concepts of moral character may serve as signal of moral
¹¹² reputation (for others) or moral self-concept. Indeed, previous studies used the moral
¹¹³ words and found that moral related information can be perceived faster (Gantman & Van
¹¹⁴ Bavel, 2014, but see, @firestone_enhanced_2015). If moral character is an important in
¹¹⁵ person perception, then, just as those other information such as races and stereotype (see

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

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

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

¹⁴² **Morality as a social-categorization?**

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

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

¹⁶³ However, the above four experiments can not distinguish between these two
¹⁶⁴ possibilities, because the concept “good-peron” can both be rewarding and be categorized
¹⁶⁵ as ingroup memeber, and previous studies using associative learning paradigm revealed
¹⁶⁶ that both rewarding stimuli (e.g., Sui, He, & Humphreys, 2012) and in-group information
¹⁶⁷ [Enock et al. (2018); enock_overlap_2020] are prioritized.

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

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

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

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

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

222 Key concepts and discussing points:

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

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

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

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

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

233 Morality as a way for social-categorization (McHugh, McGann, Igou, & Kinsella,
234 2019)? People are more likely to identify themselves with trustworthy faces (Verosky &
235 Todorov, 2010) (trustworthy faces has longer RTs).

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

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

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

242 Disclosures

243 We reported all the measurements, analyses, and results in all the experiments in the
244 current study. Participants whose overall accuracy lower than 60% were excluded from

245 analysis. Also, the accurate responses with less than 200ms reaction times were excluded
246 from the analysis.

247 All the experiments reported were not pre-registered. Most experiments (1a ~ 6b,
248 except experiment 3b) reported in the current study were first finished between 2014 to
249 2016 in Tsinghua University, Beijing, China. Participants in these experiments were
250 recruited in the local community. To increase the sample size of experiments to 50 or more
251 (Simmons, Nelson, & Simonsohn, 2013), we recruited additional participants in Wenzhou
252 University, Wenzhou, China in 2017 for experiment 1a, 1b, 4a, and 4b. Experiment 3b was
253 finished in Wenzhou University in 2017. To have a better estimation of the effect size, we
254 included the data from two experiments (experiment 7a, 7b) that were reported in Hu et
255 al. (2020) (See Table S1 for overview of these experiments).

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

258 **General methods**

259 **Design and Procedure**

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

270 higher accuracy at the end of the practicing session, they started the experimental task
271 which was the same as in the practice phase.

272 The experiment 1a, 1b, 1c, 2, and 6a shared a 2 (matching: match vs. nonmatch) by
273 3 (moral valence: good vs. neutral vs. bad) within-subject design. Experiment 1a was the
274 first one of the whole series studies and 1b, 1c, and 2 were conducted to exclude the
275 potential confounding factors. More specifically, experiment 1b used different Chinese
276 words as label to test whether the effect only occurred with certain familiar words.
277 Experiment 1c manipulated the moral valence indirectly: participants first learned to
278 associate different moral behaviors with different neutral names, after remembered the
279 association, they then performed the perceptual matching task by associating names with
280 different shapes. Experiment 2 further tested whether the way we presented the stimuli
281 influence the effect of valence, by sequentially presenting labels and shapes. Note that part
282 of participants of experiment 2 were from experiment 1a because we originally planned a
283 cross task comparison. Experiment 6a, which shared the same design as experiment 2, was
284 an EEG experiment which aimed at exploring the neural correlates of the effect. But we
285 will focus on the behavioral results of experiment 6a in the current manuscript.

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

297 3b was designed to separate the self-referential trials and other-referential trials. That is,
298 participants finished two different blocks: in the self-referential blocks, they only responded
299 to good-self, neutral-self, and bad-self, with half match trials and half non-match trials; for
300 the other-reference blocks, they only responded to good-other, neutral-other, and
301 bad-other. Experiment 7a and 7b were designed to test the cross task robustness of the
302 effect we observed in the aforementioned experiments (see, Hu et al., 2020). The matching
303 task in these two experiments shared the same design with experiment 3a, but only with
304 two moral valence, i.e., good vs. bad. We didn't include the neutral condition in
305 experiment 7a and 7b because we found that the neutral and bad conditions constantly
306 showed non-significant results in experiment 1 ~ 6.

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

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

321 E-prime 2.0 was used for presenting stimuli and collecting behavioral responses,
322 except that experiment 7a and 7b used Matlab Psychtoolbox (Brainard, 1997; Pelli, 1997).

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

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

340 Data analysis

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

345 *Classic NHST.*

346 First, as in Sui et al. (2012), we analyzed the accuracy and reaction times using
347 classic repeated measures ANOVA in the Null Hypothesis Significance Test (NHST)

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

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

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

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

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

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

373 The classic NHST approach may ignore the uncertainty in estimate of the parameters
 374 for SDT (Rouder & Lu, 2005), and using mean RT assumes normal distribution of RT
 375 data, which is always not true because RTs distribution is skewed (Rousselet & Wilcox,
 376 2019). To better estimate the uncertainty and use a more appropriate model, we also tried
 377 Bayesian hierarchical generalized linear model to analyze each experiment's accuracy and
 378 RTs data. We used BRMs (Bürkner, 2017) to build the model, which used Stan (Carpenter
 379 et al., 2017) to estimate the posterior.

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

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

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

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

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

392 In each experiment, we had multiple participants, then we need also consider the
 393 variations between subjects, i.e., a hierarchical mode in which individual's parameter and
 394 the population level parameter are estimated simultaneously. We assume that the

³⁹⁵ outcome of each trial is Bernoulli distributed (binomial with 1 trial), with probability p_{ij}
³⁹⁶ that $y_{ij} = 1$.

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

³⁹⁷ Similarly, the generalized linear model was extended to two levels:

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

³⁹⁸ The outcomes y_{ij} are 0 if participant j responded “nonmatch” on trial i , 1 if they
³⁹⁹ responded “match”. The probability of the “match” response for trial i for subject j is p_{ij} .
⁴⁰⁰ We again can write the generalized linear model on the probits (z-scores; Φ , “Phi”) of ps .

⁴⁰¹ The subjective-specific intercepts ($\beta_0 = -zFAR$) and slopes ($\beta_1 = d'$) are described
⁴⁰² by multivariate normal with means and a covariance matrix for the parameters.

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

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)). This distribution has
⁴⁰⁵ two parameters: μ , σ . μ is the mean of the logNormal distribution, and σ is the disperse of
⁴⁰⁶ the distribution. The log normal distribution can be extended to shifted log normal
⁴⁰⁷ distribution, with one more parameter: shift, which is the earliest possible response.

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

⁴⁰⁸ Shifted log-normal distribution:

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

⁴⁰⁹ y_{ij} is the RT of the i th trial of the j th participants.

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

$$\sigma_j \sim Cauchy()$$

410 *Hierarchical drift diffusion model (HDDM).*

411 To further explore the psychological mechanism under perceptual decision-making, we
 412 used HDDM (Wiecki, Sofer, & Frank, 2013) to model our RTs and accuracy data. We used
 413 the prior implemented in HDDM, that is, informative priors that constrains parameter
 414 estimates to be in the range of plausible values based on past literature (Matzke &
 415 Wagenmakers, 2009). As reported in Hu et al. (2020), we used the response code approach,
 416 match response were coded as 1 and nonmatch responses were coded as 0. To fully explore
 417 all parameters, we allow all four parameters of DDM free to vary. We then extracted the
 418 estimation of all the four parameters for each participants for the correlation analyses.

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

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

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

⁴³² self-relevance. Cohen's d and its variance were estimated using the following formula

⁴³³ (Cooper, Hedges, & Valentine, 2009):

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

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

⁴³⁴ M_1 is the mean of the first condition, sd_1 is the standard deviation of the first
⁴³⁵ condition, while M_2 is the mean of the second condition, sd_2 is the standard deviation of
⁴³⁶ the second condition. r is the correlation coefficient between data from first and second
⁴³⁷ condition. n is the number of data point (in our case the number of participants included
⁴³⁸ in our research).

⁴³⁹ The effect size from each experiment were then synthesized by random effect model
⁴⁴⁰ using `metafor` (Viechtbauer, 2010). Note that to avoid the cases that some participants
⁴⁴¹ participated more than one experiments, we inspected the all available information of
⁴⁴² participants and only included participants' results from their first participation. As
⁴⁴³ mentioned above, 24 participants were intentionally recruited to participate both exp 1a
⁴⁴⁴ and exp 2, we only included their results from experiment 1a in the meta-analysis.

⁴⁴⁵ We also estimated the synthesized effect size using Bayesian hierarchical model,
⁴⁴⁶ which extended the two-level hierarchical model in each experiment into three-level model,
⁴⁴⁷ which experiment as an additional level. For SDT, we can use a nested hierarchical model
⁴⁴⁸ to model all the experiment with similar design:

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

⁴⁴⁹ where

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

- ₄₅₀ The outcomes y_{ijk} are 0 if participant j in experiment k responded “nonmatch” on trial i ,
₄₅₁ 1 if they responded “match”.

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

- ₄₅₂ and the experiment level parameter mu_{0k} and mu_{1k} is from a higher order
₄₅₃ distribution:

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

- ₄₅₄ in which μ_0 and μ_1 means the population level parameter.

- ₄₅₅ This model can be easily expand to three-level model in which participants and
₄₅₆ experiments are two group level variable and participants were nested in the experiments.

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

- ₄₅₇ y_{ijk} is the RT of the i th trial of the j th participants in the k th experiment.

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

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

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

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

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

462 ***Valence effect.***

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

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

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

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

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

482 ***Specificity of the valence effect.***

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

487 ***Behavior-Questionnaire correlation.***

488 Finally, we explored correlation between results from behavioral results and
489 self-reported measures.

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

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

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

506 Then we confirmed the significant results from the data without self- and
507 other-referential (exp1a, 1b, 1c, 2, 5, 6a). In this confirmatory analysis, we used $\alpha =$
508 0.05 and used bootstrap by BootES package (Kirby & Gerlanc, 2013) to estimate the

correlation. To avoid false positive, we further determined the threshold for significant by permutation. More specifically, for each pairs that initially with $p < .05$, we randomly shuffle the participants data of each score and calculated the correlation between the shuffled vectors. After repeating this procedure for 5000 times, we choose arrange these 5000 correlation coefficients and use the 95% percentile number as our threshold.

Part 1: Moral valence effect

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

Experiment 1a

Methods.

Participants.

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

Stimuli and Tasks.

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

532 ***Procedure.***

533 This experiment had two phases. First, there was a brief learning stage. Participants

534 were asked to learn the relationship between geometric shapes (triangle, square, and circle)

535 and different person (bad person, a good person, or a neutral person). For example, a

536 participant was told, “bad person is a circle; good person is a triangle; and a neutral person

537 is represented by a square.” After participant remember the associations (usually in a few

538 minutes), participants started a practicing phase of matching task which has the exact task

539 as in the experimental task. In the experimental task, participants judged whether

540 shape–label pairs, which were subsequently presented, were correct. Each trial started with

541 the presentation of a central fixation cross for 500 ms. Subsequently, a pairing of a shape

542 and label (good person, bad person, and neutral person) was presented for 100 ms. The

543 pair presented could confirm to the verbal instruction for each pairing given in the training

544 stage, or it could be a recombination of a shape with a different label, with the shape–label

545 pairings being generated at random. The next frame showed a blank for 1100ms.

546 Participants were expected to judge whether the shape was correctly assigned to the person

547 by pressing one of the two response buttons as quickly and accurately as possible within

548 this timeframe (to encourage immediate responding). Feedback (correct or incorrect) was

549 given on the screen for 500 ms at the end of each trial, if no response detected, “too slow”

550 was presented to remind participants to accelerate. Participants were informed of their

551 overall accuracy at the end of each block. The practice phase finished and the experimental

552 task began after the overall performance of accuracy during practice phase achieved 60%.

553 For participants from the Tsinghua community, they completed 6 experimental blocks of 60

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

555 nonmatch, good-person match, good-person nonmatch, neutral-person match, and

556 neutral-person nonmatch). For the participants from Wenzhou University, they finished 6

557 blocks of 120 trials, therefore, 120 trials for each condition.

558 ***Data analysis.***

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

562 **Results.**

563 ***Classic NHST.***

564 *d prime.*

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

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

574 *Reaction times.*

575 We conducted 2 (Matchness: match v. nonmatch) by 3 (Valence: good, neutral, bad)
576 repeated measure ANOVA. We found the main effect of Matchness ($F(1, 50) = 232.39,$
577 $MSE = 948.92, p < .001, \hat{\eta}_G^2 = .104$), main effect of valence ($F(1.87, 93.31) = 9.62,$
578 $MSE = 1,673.86, p < .001, \hat{\eta}_G^2 = .016$), and interaction between Matchness and Valence
579 ($F(1.73, 86.65) = 8.52, MSE = 1,441.75, p = .001, \hat{\eta}_G^2 = .011$).

580 We then carried out two separate ANOVA for Match and Mismatched trials. For
581 matched trials, we found the effect of valence . We further examined the effect of valence
582 for both self and other for matched trials. We found that shapes associated with Good
583 Person (684 ms, SE = 11.5) responded faster than Neutral (709 ms, SE = 11.5), $t(50) =$

584 -2.265, $p = 0.0702$) and Bad Person (728 ms, SE = 11.7), $t(50) = -4.41$, $p = 0.0002$), and
 585 the Neutral condition was faster than the Bad condition, $t(50) = -2.495$, $p = 0.0415$). For
 586 non-matched trials, there was no significant effect of Valence ().

587 ***Bayesian hierarchical GLM.***

588 *d prime.*

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

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

602 *Reaction times.*

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

610 mismatched trials are largely overlapped. See Figure 2.

611 **HDDM.**

612 We fitted our data with HDDM, using the response-coding (See also, Hu et al., 2020).

613 We estimated separate drift rate (v), non-decision time (T_0), and boundary separation (a)

614 for each condition. We found that the shapes tagged with good person has higher drift rate

615 and higher boundary separation than shapes tagged with both neutral and bad person.

616 Also, the shapes tagged with neutral person has a higher drift rate than shapes tagged

617 with bad person, but not for the boundary separation. Finally, we found that shapes

618 tagged with bad person had longer non-decision time (see Figure 3).

619 **Experiment 1b**

620 In this study, we aimed at excluding the potential confounding factor of the

621 familiarity of words we used in experiment 1a, by matching the familiarity of the words.

622 **Method.**

623 **Participants.**

624 72 college students (49 female, age = 20.17 ± 2.08 years) participated. 39 of them

625 were recruited from Tsinghua University community in 2014; 33 were recruited from

626 Wenzhou University in 2017. All participants were right-handed except one, and all had

627 normal or corrected-to-normal vision. Informed consent was obtained from all participants

628 prior to the experiment according to procedures approved by the local ethics committees.

629 20 participant's data were excluded from analysis because nearly random level of accuracy,

630 leaving 52 participants (36 female, age = 20.25 ± 2.31 years).

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

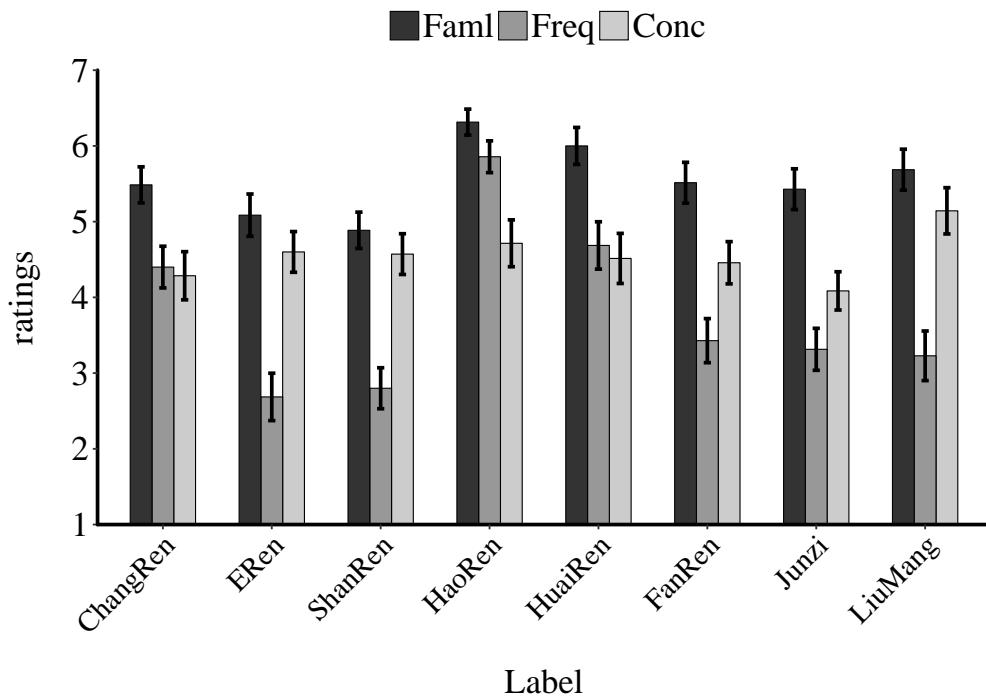
632 $\times 3.7^\circ$ of visual angle) were presented above a white fixation cross subtending $0.8^\circ \times 0.8^\circ$

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

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

635 good person (“ ”, ShanRen) or a morally neutral person (“ ”, ChangRen). The order of
 636 the associations between shapes and labels was counterbalanced across participants. Three
 637 labels used in this experiment is selected based on the rating results from an independent
 638 survey, in which participants rated the familiarity, frequency, and concreteness of eight
 639 different words online. Of the eight words, three of them are morally positive (HaoRen,
 640 ShanRen, Junzi), two of them are morally neutral (ChangRen, FanRen), and three of them
 641 are morally negative (HuaiRen, ERen, LiuMang). An independent sample consist of 35
 642 participants (22 females, age 20.6 ± 3.11) were recruited to rate these words. Based on the
 643 ratings (see supplementary materials Figure S1), we selected ShanRen, ChangRen, and
 644 ERen to represent morally positive, neutral, and negative person.

Ratings for each label



645

Procedure.

646 For participants from both Tsinghua community and Wenzhou community, the
 647 procedure in the current study was exactly same as in experiment 1a.
 648

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

650 **Results.**

651 **NHST.**

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

653 d prime.

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

655 $MSE = 0.18$, $p < .001$, $\hat{\eta}_G^2 = .053$. Paired t test showed that the Good-Person condition

656 (1.87 ± 0.102) was with greater d prime than Neutral condition $(1.44 \pm 0.101$, $t(51) =$

657 5.945 , $p < 0.001$). We also found that the Bad-Person condition (1.67 ± 0.11) has also

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

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

660 *Reaction times.*

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

662 $MSE = 960.92$, $p < .001$, $\hat{\eta}_G^2 = .031$) and then analyzed the matched trials and

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

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

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

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

667 $t(51) = -5.724$, $p < 0.0001$). While there was no significant differences between Neutral and

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

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

670 **BGLM.**

671 *Signal detection theory analysis of accuracy.*

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

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

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

675 person (2.46, 95% CI[2.21 2.72]) is greater than shapes tagged with moral bad (2.07, 95%
 676 CI[1.83 2.32]), $P_{PosteriorComparison} = 1$. Shape tagged with morally good person is also
 677 greater than shapes tagged with neutral person (2.23, 95% CI[1.95 2.49]),
 678 $P_{PosteriorComparison} = 0.97$. Also, the shapes tagged with neutral person is greater than
 679 shapes tagged with morally bad person, $P_{PosteriorComparison} = 0.92$.

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

685 *Reaction time.*

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

694 **HDDM.**

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

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

705 **Experiment 1c**

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

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

710 **Method.**

711 ***Participants.***

712 23 college students (15 female, age = 22.61 ± 2.62 years) participated. All of them
713 were recruited from Tsinghua University community in 2014. Informed consent was
714 obtained from all participants prior to the experiment according to procedures approved by
715 the local ethics committees. No participant was excluded because they overall accuracy
716 were above 0.6.

717 ***Stimuli and Tasks.***

718 Three geometric shapes (triangle, square, and circle, with $3.7^\circ \times 3.7^\circ$ of visual angle)
719 were presented above a white fixation cross subtending $0.8^\circ \times 0.8^\circ$ of visual angle at the
720 center of the screen. The three most common names were chosen, which are neutral in
721 moral valence before the manipulation. Three names (Zhang, Wang, Li) were first paired
722 with three paragraphs of behavioral description. Each description includes one sentence of
723 biographic information and four sentences that describing the moral behavioral under that
724 name. To assess the that these three descriptions represented good, neutral, and bad

valence, we collected the ratings of three person on six dimensions: morality, likability, trustworthiness, dominance, competence, and aggressiveness, from an independent sample ($n = 34$, 18 female, age = 19.6 ± 2.05). The rating results showed that the person with morally good behavioral description has higher score on morality ($M = 3.59$, $SD = 0.66$) than neutral ($M = 0.88$, $SD = 1.1$), $t(33) = 12.94$, $p < .001$, and bad conditions ($M = -3.4$, $SD = 1.1$), $t(33) = 30.78$, $p < .001$. Neutral condition was also significant higher than bad conditions $t(33) = 13.9$, $p < .001$ (See supplementary materials).

Procedure.

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

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

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

⁷⁵¹ $MSE = 1,973.91$, $p = .308$, $\hat{\eta}_G^2 = .005$).

⁷⁵² ***Signal detection theory analysis of accuracy.***

⁷⁵³ We fitted a Bayesian hierarchical GLM for SDT. The results showed that when the
⁷⁵⁴ shapes were tagged with labels with different moral valence, the sensitivity (d') and criteria
⁷⁵⁵ (c) were both influenced. For the d' , we found that the shapes tagged with morally good
⁷⁵⁶ person (2.30, 95% CI[1.93 2.70]) is greater than shapes tagged with moral bad (2.11, 95%
⁷⁵⁷ CI[1.83 2.42]), $P_{PosteriorComparison} = 0.8$. Shape tagged with morally good person is also
⁷⁵⁸ greater than shapes tagged with neutral person (2.16, 95% CI[1.88 2.45]),
⁷⁵⁹ $P_{PosteriorComparison} = 0.75$.

⁷⁶⁰ Interesting, we also found the criteria for three conditions also differ, the shapes
⁷⁶¹ tagged with good person has the highest criteria (-0.97, [-1.12 -0.82]), followed by shapes
⁷⁶² tagged with neutral person(-0.96, [-1.09 -0.83]), and then the shapes tagged with bad
⁷⁶³ person(-1.03, [-1.22 -0.84]). However, pair-wise comparison showed that only showed strong
⁷⁶⁴ evidence for the difference between good and bad conditions.

⁷⁶⁵ ***Reaction time.***

⁷⁶⁶ We fitted a Bayesian hierarchical GLM for RTs, with a log-normal distribution as the
⁷⁶⁷ link function. We used the posterior distribution of the regression coefficient to make
⁷⁶⁸ statistical inferences. As in previous studies, the matched conditions are much faster than
⁷⁶⁹ the mismatched trials ($P_{PosteriorComparison} = .75$). We focused on matched trials only, and
⁷⁷⁰ compared different conditions: Good () is not faster than the neutral (),
⁷⁷¹ $P_{PosteriorComparison} = .5$, it was faster than the Bad condition (),
⁷⁷² $P_{PosteriorComparison} = .88$. And the neutral condition is faster than the bad condition,
⁷⁷³ $P_{PosteriorComparison} = .95$. However, the mismatched trials are largely overlapped.

⁷⁷⁴ **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et
⁷⁷⁵ al., 2020). We estimated separate drift rate (v), non-decision time (T_0), and boundary
⁷⁷⁶ separation (a) for each condition. We found that the shapes tagged with good person has

777 higher drift rate and higher boundary separation than shapes tagged with both neutral and
778 bad person. Also, the shapes tagged with neutral person has a higher drift rate than
779 shapes tagged with bad person, but not for the boundary separation. Finally, we found
780 that shapes tagged with bad person had longer non-decision time (see figure 9)).

781 **Experiment 2: Sequential presenting**

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

794 **Method.**

795 ***Participants.***

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

801 ***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 2×2 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.

815 Data analysis.

816 Data was analyzed as in study 1a.

817 Results.

818 NHST.

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

821 d prime.

822 There was evidence for the main effect of valence, $F(1.83, 60.36) = 14.41$,
823 $MSE = 0.23$, $p < .001$, $\eta^2_G = .066$. Paired t test showed that the Good-Person condition
824 (2.79 ± 0.17) was with greater d prime than Netural condition (2.21 ± 0.16 , $t(33) = 4.723$,
825 $p = 0.001$) and Bad-person condition (2.41 ± 0.14), $t(33) = 4.067$, $p = 0.008$). There was
826 no-significant difference between Neutral-person and Bad-person condidition, $t(33) = -1.802$,
827 $p = 0.185$.

828 *Reaction time.*

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

831 We found interaction between Matchness and Valence ($F(1.99, 65.70) = 9.53$,
 832 $MSE = 605.36$, $p < .001$, $\hat{\eta}_G^2 = .017$) and then analyzed the matched trials and
 833 mismatched trials separately, as in experiment 1a. For matched trials, we found the effect
 834 of valence $F(1.99, 65.76) = 10.57$, $MSE = 1,192.65$, $p < .001$, $\hat{\eta}_G^2 = .067$. Post-hoc t -tests
 835 revealed that shapes associated with Good Person (548 ± 9.4) were responded faster than
 836 Neutral-Person (582 ± 10.9), ($t(33) = -3.95$, $p = 0.0011$) and Bad Person (582 ± 10.2),
 837 $t(33) = -3.9$, $p = 0.0013$). While there was no significant differences between Neutral and
 838 Bad-Person condition ($t(33) = -0.01$, $p = 0.999$). For non-matched trials, there was no
 839 significant effect of Valence ($F(1.99, 65.83) = 0.17$, $MSE = 489.80$, $p = .843$, $\hat{\eta}_G^2 = .001$).

840 **BGLMM.**

841 *Signal detection theory analysis of accuracy.*

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

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

854 evidence for the difference between good and bad conditions.

855 *Reaction times.*

856 We fitted a Bayesian hierarchical GLM for RTs, with a log-normal distribution as the
857 link function. We used the posterior distribution of the regression coefficient to make
858 statistical inferences. As in previous studies, the matched conditions are much faster than
859 the mismatched trials ($P_{PosteriorComparison} = .75$). We focused on matched trials only, and
860 compared different conditions: Good () is not faster than the neutral (),
861 $P_{PosteriorComparison} = .5$, it was faster than the Bad condition (),
862 $P_{PosteriorComparison} = .88$. And the neutral condition is faster than the bad condition,
863 $P_{PosteriorComparison} = .95$. However, the mismatched trials are largely overlapped.

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

872 Discussion

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

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

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

882 **Method.**

883 ***Participants.***

884 24 college students (8 female, age = 22.88 ± 2.79) participated the current study, all
885 of them were from Tsinghua University in 2014. Informed consent was obtained from all
886 participants prior to the experiment according to procedures approved by a local ethics
887 committee. No participant was excluded from behavioral analysis.

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

892 ***Stimuli.***

893 Three geometric shapes (triangle, square and circle, each $4.6^\circ \times 4.6^\circ$ of visual angle)
894 were presented at the center of screen for 50 ms after 500ms of fixation ($0.8^\circ \times 0.8^\circ$ of
895 visual angle). The association of the three shapes to bad person (“ , HuaiRen”), good
896 person (“ , HaoRen”) or ordinary person (“ , ChangRen”) was counterbalanced across
897 participants. The words bad person, good person or ordinary person ($3.6^\circ \times 1.6^\circ$) was also
898 displayed at the center fo the screen. Participants had to judge whether the pairings of
899 label and shape matched (e.g., Does the circle represent a bad person?). The experiment
900 was run on a PC using E-prime software (version 2.0). These stimuli were displayed on a
901 22-in CRT monitor (1024×768 at 100Hz). We used backward masking to avoid
902 over-processing of the moral words, in which a scrambled picture were presented for 900 ms
903 after the label. Also, to avoid the ceiling effect on accuracy, shapes were presented on a

904 noisy background based on our pilot studies. The noisy images were made by scrambling a
905 picture of 3/4 gray and 1/4 white at resolution of 2×2 pixel.

906 ***Procedure.***

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

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

918 All the stimuli were presented on a gray background (RGB: 127, 127, 127). E-primed
919 2.0 was used to present stimuli and collect behavioral results. Data were collected and
920 analyzed when accuracy performance in total reached 60%.

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

922 **Results.**

923 **NHST.**

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

926 d prime.

927 We conducted repeated measures ANOVA, with moral valence as independent
928 variable. The results revealed the main effect of valence ($F(1.74, 40.05) = 3.76$,

929 $MSE = 0.10, p = .037, \hat{\eta}_G^2 = .021$). Post-hoc analysis revealed that shapes link with Good
 930 person (mean = 3.13, SE = 0.109) is greater than Neutral condition (mean = 2.88, SE =
 931 0.14), $t = 2.916, df = 24, p = 0.02$, p-value adjusted by Tukey method, but the d prime
 932 between Good and bad (mean = 3.03, SE = 0.142) ($t = 1.512, df = 24, p = 0.3034$, p-value
 933 adjusted by Tukey method), bad and neutral ($t = 1.599, df = 24, p = 0.2655$, p-value
 934 adjusted by Tukey method) were not significant.

935 *Reaction times.*

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

938 We found intercation between Matchness and Valence ($F(1.97, 45.20) = 20.45$,
 939 $MSE = 450.47, p < .001, \hat{\eta}_G^2 = .021$) and then analyzed the matched trials and
 940 mismatched trials separately, as in experiment 2. For matched trials, we found the effect of
 941 valence $F(1.97, 45.25) = 32.37, MSE = 522.42, p < .001, \hat{\eta}_G^2 = .078$. For non-matched
 942 trials, there was no significant effect of Valence ($F(1.77, 40.67) = 0.35, MSE = 242.15$,
 943 $p = .679, \hat{\eta}_G^2 = .000$). Post-hoc t -tests revealed that shapes associated with Good Person
 944 (mean = 550, SE = 13.8) were responded faster than Neutral-Person (501, SE = 14.7),
 945 ($t(24) = -5.171, p = 0.0001$) and Bad Person (523, SE = 16.3), $t(24) = -8.137, p <$
 946 0.0001., and Neutral is faster than Bad-Person condition ($t(32) = -3.282, p = 0.0085$).

947 **BGLM.**

948 *Signal detection theory analysis of accuracy.*

949 *Reaction time.*

950 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et
 951 al., 2020). We estimated separate drift rate (v), non-decision time (T_0), and boundary
 952 separation (a) for each condition. We found that, similar to experiment 2, the shapes
 953 tagged with good person has higher drift rate and higher boundary separation than shapes
 954 tagged with both neutral and bad person, but only for the self-referential condition. Also,

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

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

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

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

966 **Experiment 3a**

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

971 **Method.**

972 ***Participants.***

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

Design.

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

Stimuli.

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

Procedure.

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

¹⁰⁰⁵ non-matched trials for each condition (good self, neutral self, bad self, good other, neutral
¹⁰⁰⁶ other, bad other) for each block.

¹⁰⁰⁷ ***Data Analysis.***

¹⁰⁰⁸ Data analysis followed strategies described in the general method section. Reaction
¹⁰⁰⁹ times and d prime data were analyzed as in study 1 and study 2, except that one more
¹⁰¹⁰ within-subject variable (i.e., self-relevance) was included in the analysis.

¹⁰¹¹ **Results.**

¹⁰¹² ***NHST.***

¹⁰¹³ Figure 16 shows d prime and reaction times of experiment 3a. Less than 5% correct
¹⁰¹⁴ trials with less than 200ms reaction times were excluded.

¹⁰¹⁵ *d prime.*

¹⁰¹⁶ There was evidence for the main effect of valence, $F(1.89, 64.37) = 11.09$,
¹⁰¹⁷ $MSE = 0.23$, $p < .001$, $\hat{\eta}_G^2 = .039$, and main effect of self-relevance, $F(1, 34) = 3.22$,
¹⁰¹⁸ $MSE = 0.54$, $p = .082$, $\hat{\eta}_G^2 = .015$, as well as the interaction, $F(1.79, 60.79) = 3.39$,
¹⁰¹⁹ $MSE = 0.43$, $p = .045$, $\hat{\eta}_G^2 = .022$.

¹⁰²⁰ We then conducted separated ANOVA for self-referential and other-referential trials.
¹⁰²¹ The valence effect was shown for the self-referential conditions, $F(1.65, 56.25) = 13.98$,
¹⁰²² $MSE = 0.31$, $p < .001$, $\hat{\eta}_G^2 = .119$. Post-hoc test revealed that the Good-Self condition
¹⁰²³ (1.97 ± 0.14) was with greater d prime than Netural condition (1.41 ± 0.12 , $t(34) = 4.505$,
¹⁰²⁴ $p = 0.0002$), and Bad-self condition (1.43 ± 0.102), $t(34) = 3.856$, $p = 0.0014$. There was
¹⁰²⁵ difference between neutral and bad condition, $t(34) = -0.238$, $p = 0.9694$. However, no
¹⁰²⁶ effect of valence was found for the other-referential condition $F(1.98, 67.36) = 0.38$,
¹⁰²⁷ $MSE = 0.35$, $p = .681$, $\hat{\eta}_G^2 = .004$.

¹⁰²⁸ *Reaction time.*

¹⁰²⁹ We found interaction between Matchness and Valence ($F(1.98, 67.44) = 26.29$,

₁₀₃₀ $MSE = 730.09, p < .001, \hat{\eta}_G^2 = .025$) and then analyzed the matched trials and nonmatch
₁₀₃₁ trials separately, as in previous experiments.

₁₀₃₂ For the match trials, we found that the interaction between identity and valence,
₁₀₃₃ $F(1.72, 58.61) = 3.89, MSE = 2,750.19, p = .032, \hat{\eta}_G^2 = .019$, as well as the main effect of
₁₀₃₄ valence $F(1.98, 67.34) = 35.76, MSE = 1,127.25, p < .001, \hat{\eta}_G^2 = .079$, but not the effect of
₁₀₃₅ identity $F(1, 34) = 0.20, MSE = 3,507.14, p = .660, \hat{\eta}_G^2 = .001$. As for the d prime, we
₁₀₃₆ separated analyzed the self-referential and other-referential trials. For the Self-referential
₁₀₃₇ trials, we found the main effect of valence, $F(1.80, 61.09) = 30.39, MSE = 1,584.53,$
₁₀₃₈ $p < .001, \hat{\eta}_G^2 = .159$; for the other-referential trials, the effect of valence is weaker,
₁₀₃₉ $F(1.86, 63.08) = 2.85, MSE = 2,224.30, p = .069, \hat{\eta}_G^2 = .024$. We then focused on the self
₁₀₄₀ conditions: the good-self condition (713 ± 12) is faster than neutral- (776 ± 11.8), $t(34) =$
₁₀₄₁ $-7.396, p < .0001$, and bad-self (772 ± 10.1) conditions, $t(34) = -5.66, p < .0001$. But
₁₀₄₂ there is not difference between neutral- and bad-self conditions, $t(34) = 0.481, p = 0.881$.

₁₀₄₃ For the nonmatch trials, we didn't found any strong effect: identity, $F(1, 34) = 3.43,$
₁₀₄₄ $MSE = 660.02, p = .073, \hat{\eta}_G^2 = .004$, valence $F(1.89, 64.33) = 0.40, MSE = 444.10,$
₁₀₄₅ $p = .661, \hat{\eta}_G^2 = .001$, or interaction between the two $F(1.94, 66.02) = 2.42, MSE = 817.35,$
₁₀₄₆ $p = .099, \hat{\eta}_G^2 = .007$.

₁₀₄₇ **BGLM.**

₁₀₄₈ *Signal detection theory analysis of accuracy.*

₁₀₄₉ We found that the d prime is greater when shapes were associated with good self
₁₀₅₀ condition than with neutral self or bad self, but shapes associated with bad self and neutral
₁₀₅₁ self didn't show differences. Comparing the self vs other under three condition revealed
₁₀₅₂ that shapes associated with good self is greater than with good other, but with a weak
₁₀₅₃ evidence. In contrast, for both neutral and bad valence condition, shapes associated with
₁₀₅₄ other had greater d prime than with self.

₁₀₅₅ *Reaction time.*

1056 In reaction times, we found that same trends in the match trials as in the RT: while
1057 the shapes associated with good self was greater than with good other (log mean diff =
1058 -0.02858, 95%HPD[-0.070898, 0.0154]), the direction is reversed for neutral and negative
1059 condition. see Figure 17

1060 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et
1061 al., 2020). We estimated separate drift rate (v), non-decision time (T_0), and boundary
1062 separation (a) for each condition. We found that the shapes tagged with good person has
1063 higher drift rate and higher boundary separation than shapes tagged with both neutral and
1064 bad person. Also, the shapes tagged with neutral person has a higher drift rate than
1065 shapes tagged with bad person, but not for the boundary separation. Finally, we found
1066 that shapes tagged with bad person had longer non-decision time (see figure 18)).

1067 **Experiment 3b**

1068 In study 3a, participants had to remember 6 pairs of association, which cause high
1069 cognitive load during the whole experiment. To eliminate the influence of cognitive load, we
1070 conducted study 3b, in which participant learn three aspect of self and stranger separately
1071 in to consecutive task. We hypothesize that we will replicate the pattern of study 3a, i.e.,
1072 the effect of moral valence only occurs for self-relevant conditions. ### Method

1073 **Participants.**

1074 Study 3b were finished in 2017, at that time we have calculated that the effect size
1075 (Cohen's d) of good-person (or good-self) vs. bad-person (or bad-other) was between 0.47 ~
1076 0.53, based on the data from Tsinghua community in study 1a, 1b, 2, 3a, 4a, and 4b. Based
1077 on this effect size, we estimated that 54 participants would allow we to detect the effect
1078 size of Cohen's $= 0.5$ with 95% power and alpha = 0.05, using G*power 3.192 (Faul,
1079 Erdfelder, Buchner, & Lang, 2009). Therefore, we planned to stop after we arrived this
1080 number. During the data collected at Wenzhou University, 61 participants (45 females; 19

1081 to 25 years of age, age = 20.42 ± 1.77) came to the testing room and we tested all of them
1082 during a single day. All participants were right-handed, and all had normal or
1083 corrected-to-normal vision. Informed consent was obtained from all participants prior to
1084 the experiment according to procedures approved by a local ethics committee. 4
1085 participants' data were excluded from analysis because their over all accuracy was lower
1086 than 60%, 1 more participant was excluded because of zero hit rate for one condition,
1087 leaving 56 participants (43 females; 19 to 25 years old, age = 20.27 ± 1.60).

1088 ***Design.***

1089 Study 3b has the same experimental design as 3a, with a $2 \times 3 \times 2$ within-subject
1090 design. The first variable was self-relevance, include two levels: self-relevant
1091 vs. stranger-relevant; the second variable was moral valence, include good, neutral and bad;
1092 the third variable was the matching between shape and label: match vs. mismatch.
1093 Stimuli. The stimuli used in study 3b share the same parameters with experiment 3a. 6
1094 shapes were included (triangle, square, circle, trapezoid, diamond, regular pentagon), as
1095 well as 6 labels, but the labels changed to "good self", "neutral self", "bad self", "good
1096 him/her", "bad him/her", "neutral him/her", the stranger's label is consistent with
1097 participants' gender. Same as study 3a, we asked participant to chosen an unfamiliar name
1098 of their own gender to be the stranger before showing them the relationship. Note, because
1099 of implementing error, the personal distance data did not collect for this experiment.

1100 ***Stimuli.***

1101 The stimuli used in study 3b is the same as in experiment 3a.

1102 ***Procedure.***

1103 In this experiment, participants finished two matching tasks, i.e., self-matching task,
1104 and other-matching task. In the self-matching task, participants first associate the three
1105 aspects of self to three different shapes, and then perform the matching task. In the
1106 other-matching task, participants first associate the three aspects of the stranger to three

1107 different shapes, and then perform the matching task. The order of self-task and other-task
1108 are counter-balanced among participants. Different from experiment 3a, after presenting
1109 the stimuli pair for 100ms, participant has 1900 ms to response, and they feedback with
1110 both accuracy and reaction time. As in study 3a, before each task, the instruction showed
1111 the meaning of each label to participants. The self-matching task and other-matching task
1112 were randomized between participants. Each participant finished 6 blocks, each have 120
1113 trials.

1114 ***Data Analysis.***

1115 Same as experiment 3a.

1116 **Results.**

1117 ***NHST.***

1118 Figure 19 shows *d* prime and reaction times of experiment 3b. Less than 5% correct
1119 trials with less than 200ms reaction times were excluded.

1120 *d prime.*

1121 There was no evidence for the main effect of valence, $F(1.92, 105.43) = 1.90$,
1122 $MSE = 0.33$, $p = .157$, $\hat{\eta}_G^2 = .005$, but we found a main effect of self-relevance,
1123 $F(1, 55) = 4.65$, $MSE = 0.89$, $p = .035$, $\hat{\eta}_G^2 = .017$, as well as the interaction,
1124 $F(1.90, 104.36) = 5.58$, $MSE = 0.26$, $p = .006$, $\hat{\eta}_G^2 = .011$.

1125 We then conducted separated ANOVA for self-referential and other-referential trials.
1126 The valence effect was shown for the self-referential conditions, $F(1.75, 96.42) = 6.73$,
1127 $MSE = 0.30$, $p = .003$, $\hat{\eta}_G^2 = .037$. Post-hoc test revealed that the Good-Self condition
1128 (2.15 ± 0.12) was with greater *d* prime than Neutral condition (1.83 ± 0.12 , $t(34) = 3.36$,
1129 $p = 0.0031$), and Bad-self condition (1.87 ± 0.12), $t(34) = 2.955$, $p = 0.01$. There was
1130 difference between neutral and bad condition, $t(34) = -0.039$, $p = 0.914$. However, no effect
1131 of valence was found for the other-referential condition $F(1.93, 105.97) = 0.61$,
1132 $MSE = 0.31$, $p = .539$, $\hat{\eta}_G^2 = .002$.

₁₁₃₃ *Reaction time.*

₁₁₃₄ We found interaction between Matchness and Valence ($F(1.86, 102.47) = 15.44$,

₁₁₃₅ $MSE = 3, 112.78, p < .001, \hat{\eta}_G^2 = .006$) and then analyzed the matched trials and

₁₁₃₆ nonmatch trials separately, as in previous experiments.

₁₁₃₇ For the match trials, we found that the interaction between identity and valence,

₁₁₃₈ $F(1.67, 92.11) = 6.14, MSE = 6, 472.48, p = .005, \hat{\eta}_G^2 = .009$, as well as the main effect of

₁₁₃₉ valence $F(1.88, 103.65) = 24.25, MSE = 5, 994.25, p < .001, \hat{\eta}_G^2 = .038$, but not the effect

₁₁₄₀ of identity $F(1, 55) = 48.49, MSE = 25, 892.59, p < .001, \hat{\eta}_G^2 = .153$. As for the d prime,

₁₁₄₁ we separated analyzed the self-referential and other-referential trials. For the

₁₁₄₂ Self-referential trials, we found the main effect of valence, $F(1.66, 91.38) = 23.98$,

₁₁₄₃ $MSE = 6, 965.61, p < .001, \hat{\eta}_G^2 = .100$; for the other-referential trials, the effect of valence

₁₁₄₄ is weaker, $F(1.89, 103.94) = 5.96, MSE = 5, 589.90, p = .004, \hat{\eta}_G^2 = .014$. We then focused

₁₁₄₅ on the self conditions: the good-self condition (713 ± 12) is faster than neutral- ($776 \pm$

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

₁₁₄₇ $.0001$. But there is not difference between neutral- and bad-self conditions, $t(34) = 0.481, p$

₁₁₄₈ $= 0.881$.

₁₁₄₉ For the nonmatch trials, we didn't found any strong effect: identity, $F(1, 55) = 10.31$,

₁₁₅₀ $MSE = 24, 590.52, p = .002, \hat{\eta}_G^2 = .035$, valence $F(1.98, 108.63) = 20.57, MSE = 2, 847.51$,

₁₁₅₁ $p < .001, \hat{\eta}_G^2 = .016$, or interaction between the two $F(1.93, 106.25) = 35.51$,

₁₁₅₂ $MSE = 1, 939.88, p < .001, \hat{\eta}_G^2 = .019$.

₁₁₅₃ **BGLM.**

₁₁₅₄ *Signal detection theory analysis of accuracy.*

₁₁₅₅ We found that the d prime is greater when shapes were associated with good self

₁₁₅₆ condition than with neutral self or bad self, but shapes associated with bad self and neutral

₁₁₅₇ self didn't show differences. comparing the self vs other under three condition revealed that

₁₁₅₈ shapes associated with good self is greater than with good other, but with a weak evidence.

₁₁₅₉ In contrast, for both neutral and bad valence condition, shapes associated with other had
₁₁₆₀ greater d' prime than with self.

₁₁₆₁ *Reaction time.*

₁₁₆₂ In reaction times, we found that same trends in the match trials as in the RT: while
₁₁₆₃ the shapes associated with good self was greater than with good other (log mean diff =
₁₁₆₄ -0.02858, 95%HPD[-0.070898, 0.0154]), the direction is reversed for neutral and negative
₁₁₆₅ condition. see Figure 20

₁₁₆₆ **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et
₁₁₆₇ al., 2020). We estimated separate drift rate (v), non-decision time (T_0), and boundary
₁₁₆₈ separation (a) for each condition. We found that, similar to experiment 3a, the shapes
₁₁₆₉ tagged with good person has higher drift rate and higher boundary separation than shapes
₁₁₇₀ tagged with both neutral and bad person, but only for the self-referential condition. Also,
₁₁₇₁ the shapes tagged with neutral person has a higher drift rate than shapes tagged with bad
₁₁₇₂ person, but not for the boundary separation, and this effect also exist only for the
₁₁₇₃ self-referential condition.

₁₁₇₄ Interestingly, we found that in both self-referential and other-referential conditions,
₁₁₇₅ the shapes associated bad valence have higher drift rate and higher boundary separation.
₁₁₇₆ which might suggest that the shape associated with bad stimuli might be prioritized in the
₁₁₇₇ non-match trials (see figure 21)).

₁₁₇₈ Experiment 6b

₁₁₇₉ Experiment 6b was conducted to study the neural correlates of the prioritization
₁₁₈₀ effect of positive self, i.e., the neural underlying of the behavioral effect found int
₁₁₈₁ experiment 3a. However, as in experiment 6a, the procedure of this experiment was
₁₁₈₂ modified to adopted to ERP experiment.

₁₁₈₃ Method.

Participants.

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

Design.

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

Stimuli.

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

Procedure.

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

1210 All the stimuli were presented on a gray background (RGB: 127, 127, 127). E-primed
1211 2.0 was used to present stimuli and collect behavioral results. Data were collected and
1212 analyzed when accuracy performance in total reached 60%.

1213 Because learning 6 associations was more difficult than 3 associations and participant
1214 might have low accuracy (see experiment 3a), the current study had extended to a two-day
1215 paradigm to maximizing the accurate trials that can be used in EEG data. At the first day,
1216 participants learnt the associations and finished 9 blocks of the matching task, each had
1217 120 trials, without EEG recording. That is, each condition has 90 trials.

1218 Participants came back to lab at the second day and finish the same task again, with
1219 EEG recorded. Before the EEG experiment, each participant finished a practice session
1220 again, if their accuracy is equal or higher than 85%, they start the experiment (one
1221 participant used lower threshold 75%). Each participant finished 18 blocks, each has 90
1222 trials. One participant finished additional 6 blocks because of high error rate at the
1223 beginning, another two participant finished addition 3 blocks because of the technique
1224 failure in recording the EEG data. To increase the number of trials that can be used for
1225 EEG data analysis, matched trials has twice number as mismatched trials, therefore, for
1226 matched trials each participants finished 180 trials for each condition, for mismatched
1227 trials, each conditions has 90 trials.

1228 ***Data Analysis.***

1229 Same as experiment 3a.

1230 **Results of Day 1.**

1231 **NHST.**

1232 Figure 22 shows d prime and reaction times of experiment 3b. Less than 5% correct
1233 trials with less than 200ms reaction times were excluded.

1234 d prime.

1235 There was no evidence for the main effect of valence, $F(1.91, 40.20) = 11.98$,
 1236 $MSE = 0.15$, $p < .001$, $\hat{\eta}_G^2 = .040$, but we found a main effect of self-relevance,
 1237 $F(1, 21) = 1.21$, $MSE = 0.20$, $p = .284$, $\hat{\eta}_G^2 = .003$, as well as the interaction,
 1238 $F(1.28, 26.90) = 12.88$, $MSE = 0.21$, $p = .001$, $\hat{\eta}_G^2 = .041$.

1239 We then conducted separated ANOVA for self-referential and other-referential trials.
 1240 The valence effect was shown for the self-referential conditions, $F(1.73, 36.42) = 29.31$,
 1241 $MSE = 0.14$, $p < .001$, $\hat{\eta}_G^2 = .147$. Post-hoc test revealed that the Good-Self condition
 1242 (2.15 ± 0.12) was with greater d prime than Neutral condition (1.83 ± 0.12 , $t(34) = 3.36$,
 1243 $p = 0.0031$), and Bad-self condition (1.87 ± 0.12), $t(34) = 2.955$, $p = 0.01$. There was
 1244 difference between neutral and bad condition, $t(34) = -0.039$, $p = 0.914$. However, no effect
 1245 of valence was found for the other-referential condition $F(1.75, 36.72) = 0.00$, $MSE = 0.18$,
 1246 $p = .999$, $\hat{\eta}_G^2 = .000$.

1247 *Reaction time.*

1248 We found interaction between Matchness and Valence ($F(1.79, 37.63) = 4.07$,
 1249 $MSE = 704.90$, $p = .029$, $\hat{\eta}_G^2 = .003$) and then analyzed the matched trials and nonmatch
 1250 trials separately, as in previous experiments.

1251 For the match trials, we found that the interaction between identity and valence,
 1252 $F(1.72, 36.16) = 4.55$, $MSE = 1,560.90$, $p = .022$, $\hat{\eta}_G^2 = .015$, as well as the main effect of
 1253 valence $F(1.93, 40.55) = 9.83$, $MSE = 1,951.84$, $p < .001$, $\hat{\eta}_G^2 = .044$, but not the effect of
 1254 identity $F(1, 21) = 4.87$, $MSE = 2,032.05$, $p = .039$, $\hat{\eta}_G^2 = .012$. As for the d prime, we
 1255 separated analyzed the self-referential and other-referential trials. For the Self-referential
 1256 trials, we found the main effect of valence, $F(1.92, 40.38) = 14.48$, $MSE = 1,647.20$,
 1257 $p < .001$, $\hat{\eta}_G^2 = .112$; for the other-referential trials, the effect of valence is weaker,
 1258 $F(1.79, 37.50) = 1.04$, $MSE = 1,842.07$, $p = .356$, $\hat{\eta}_G^2 = .008$. We then focused on the self
 1259 conditions: the good-self condition (713 ± 12) is faster than neutral- (776 ± 11.8), $t(34) =$
 1260 -7.396 , $p < .0001$, and bad-self (772 ± 10.1) conditions, $t(34) = -5.66$, $p < .0001$. But

1261 there is not difference between neutral- and bad-self conditions, $t(34) = 0.481$, $p = 0.881$.

1262 For the nonmatch trials, we didn't found any strong effect: identity, $F(1, 21) = 2.76$,

1263 $MSE = 1,718.93$, $p = .112$, $\hat{\eta}_G^2 = .006$, valence $F(1.61, 33.77) = 3.81$, $MSE = 1,532.21$,

1264 $p = .041$, $\hat{\eta}_G^2 = .012$, or interaction between the two $F(1.90, 39.97) = 2.23$, $MSE = 720.80$,

1265 $p = .123$, $\hat{\eta}_G^2 = .004$.

1266 **BGLM.**

1267 *Signal detection theory analysis of accuracy.*

1268 We found that the d prime is greater when shapes were associated with good self

1269 condition than with neutral self or bad self, but shapes associated with bad self and neutral

1270 self didn't show differences. comparing the self vs other under three condition revealed that

1271 shapes associated with good self is greater than with good other, but with a weak evidence.

1272 In contrast, for both neutral and bad valence condition, shapes associated with other had

1273 greater d prime than with self.

1274 *Reaction time.*

1275 In reaction times, we found that same trends in the match trials as in the RT: while

1276 the shapes associated with good self was greater than with good other ($\log \text{mean diff} =$

1277 -0.02858 , $95\% \text{HPD}[-0.070898, 0.0154]$), the direction is reversed for neutral and negative

1278 condition. see Figure 23

1279 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et

1280 al., 2020). We estimated separate drift rate (v), non-decision time (T_0), and boundary

1281 separation (a) for each condition. We found that, similar to experiment 3a, the shapes

1282 tagged with good person has higher drift rate and higher boundary separation than shapes

1283 tagged with both neutral and bad person, but only for the self-referential condition. Also,

1284 the shapes tagged with neutral person has a higher drift rate than shapes tagged with bad

1285 person, but not for the boundary separation, and this effect also exist only for the

1286 self-referential condition.

1287 Interestingly, we found that in both self-referential and other-referential conditions,
1288 the shapes associated bad valence have higher drift rate and higher boundary separation.
1289 which might suggest that the shape associated with bad stimuli might be prioritized in the
1290 non-match trials (see figure 24).

1291 **Part 3: Implicit binding between valence and identity**

1292 In this part, we reported two studies in which the moral valence or the self-referential
1293 processing is not task-relevant. We are interested in testing whether the task-relevance will
1294 eliminate the effect observed in previous experiment.

1295 **Experiment 4a: Morality as task-irrelevant variable**

1296 In part two (experiment 3a and 3b), participants learned the association between self
1297 and moral valence directly. In Experiment 4a, we examined whether the interaction
1298 between moral valence and identity occur even when one of the variable was irrelevant to
1299 the task. In experiment 4a, participants learnt associations between shapes and self/other
1300 labels, then made perceptual match judgments only about the self or other conditions
1301 labels and shapes (cf. Sui et al. (2012)). However, we presented labels of different moral
1302 valence in the shapes, which means that the moral valence factor become task irrelevant. If
1303 the binding between moral good and self is intrinsic and automatic, then we will observe
1304 that facilitating effect of moral good for self conditions, but not for other conditions.

1305 **Method.**

1306 ***Participants.***

1307 64 participants (37 female, age = 19.70 ± 1.22) participated the current study, 32 of
1308 them were from Tsinghua University in 2015, 32 were from Wenzhou University
1309 participated in 2017. All participants were right-handed, and all had normal or
1310 corrected-to-normal vision. Informed consent was obtained from all participants prior to

1311 the experiment according to procedures approved by a local ethics committee. The data
1312 from 5 participants from Wenzhou site were excluded from analysis because their accuracy
1313 was close to chance (< 0.6). The results for the remaining 59 participants (33 female, age
1314 = 19.78 ± 1.20) were analyzed and reported.

1315 ***Design.***

1316 As in Experiment 3, a $2 \times 3 \times 2$ within-subject design was used. The first variable was
1317 self-relevance (self and stranger associations); the second variable was moral valence (good,
1318 neutral and bad associations); the third variable was the matching between shape and label
1319 (matching vs. non-match for the personal association). However, in this the task,
1320 participants only learn the association between two geometric shapes and two labels (self
1321 and other), i.e., only self-relevance were related to the task. The moral valence
1322 manipulation was achieved by embedding the personal label of the labels in the geometric
1323 shapes, see below. For simplicity, the trials where shapes where paired with self and with a
1324 word of “good person” inside were shorted as good-self condition, similarly, the trials where
1325 shapes paired with the self and with a word of “bad person” inside were shorted as bad-self
1326 condition. Hence, we also have six conditions: good-self, neutral-self, bad-self, good-other,
1327 neutral-other, and bad-other.

1328 ***Stimuli.***

1329 2 shapes were included (circle, square) and each appeared above a central fixation
1330 cross with the personal label appearing below. However, the shapes were not empty but
1331 with a two-Chinese-character word in the middle, the word was one of three labels with
1332 different moral valence: “good person”, “bad person” and “neutral person”. Before the
1333 experiment, participants learned the self/other association, and were informed to only
1334 response to the association between shapes’ configure and the labels below the fixation, but
1335 ignore the words within shapes. Besides the behavioral experiments, participants from
1336 Tsinghua community also finished questionnaires as Experiments 3, and participants from

₁₃₃₇ Wenzhou community finished a series of questionnaire as the other experiment finished in
₁₃₃₈ Wenzhou.

₁₃₃₉ ***Procedure.***

₁₃₄₀ The procedure was similar to Experiment 1. There were 6 blocks of trial, each with
₁₃₄₁ 120 trials for 2017 data. Due to procedure error, the data collected in 2015 in Tsinghua
₁₃₄₂ community only have 60 trials for each block, i.e., 30 trials per condition.

₁₃₄₃ As in study 3a, before each task, the instruction showed the meaning of each label to
₁₃₄₄ participants. The self-matching task and other-matching task were randomized between
₁₃₄₅ participants. Each participant finished 6 blocks, each have 120 trials.

₁₃₄₆ ***Data Analysis.***

₁₃₄₇ Same as experiment 3a.

₁₃₄₈ **Results.**

₁₃₄₉ ***NHST.***

₁₃₅₀ Figure 25 shows d prime and reaction times of experiment 3a. Less than 5% correct
₁₃₅₁ trials with less than 200ms reaction times were excluded.

₁₃₅₂ d prime.

₁₃₅₃ There was no evidence for the main effect of valence, $F(1.93, 111.66) = 0.53$,
₁₃₅₄ $MSE = 0.12$, $p = .581$, $\hat{\eta}_G^2 = .000$, but we found a main effect of self-relevance,
₁₃₅₅ $F(1, 58) = 121.04$, $MSE = 0.48$, $p < .001$, $\hat{\eta}_G^2 = .189$, as well as the interaction,
₁₃₅₆ $F(1.99, 115.20) = 4.12$, $MSE = 0.14$, $p = .019$, $\hat{\eta}_G^2 = .004$.

₁₃₅₇ We then conducted separated ANOVA for self-referential and other-referential trials.
₁₃₅₈ The valence effect was shown for the self-referential conditions, $F(1.95, 112.92) = 3.01$,
₁₃₅₉ $MSE = 0.15$, $p = .055$, $\hat{\eta}_G^2 = .008$. Post-hoc test revealed that the Good-Self condition
₁₃₆₀ (2.15 ± 0.12) was with greater d prime than Neutral condition (1.83 ± 0.12 , $t(34) = 3.36$,
₁₃₆₁ $p = 0.0031$), and Bad-self condition (1.87 ± 0.12), $t(34) = 2.955$, $p = 0.01$. There was

₁₃₆₂ difference between neutral and bad condition, $t(34) = -0.039$, $p = 0.914$. However, no effect
₁₃₆₃ of valence was found for the other-referential condition $F(1.98, 114.61) = 1.75$,
₁₃₆₄ $MSE = 0.10$, $p = .179$, $\hat{\eta}_G^2 = .003$.

₁₃₆₅ *Reaction time.*

₁₃₆₆ We found interaction between Matchness and Valence ($F(1.94, 112.64) = 0.84$,
₁₃₆₇ $MSE = 465.35$, $p = .432$, $\hat{\eta}_G^2 = .000$) and then analyzed the matched trials and nonmatch
₁₃₆₈ trials separately, as in previous experiments.

₁₃₆₉ For the match trials, we found that the interaction between identity and valence,
₁₃₇₀ $F(1.90, 110.18) = 4.41$, $MSE = 465.91$, $p = .016$, $\hat{\eta}_G^2 = .003$, as well as the main effect of
₁₃₇₁ valence $F(1.98, 114.82) = 0.94$, $MSE = 606.30$, $p = .392$, $\hat{\eta}_G^2 = .001$, but not the effect of
₁₃₇₂ identity $F(1, 58) = 124.15$, $MSE = 4,037.53$, $p < .001$, $\hat{\eta}_G^2 = .257$. As for the d prime, we
₁₃₇₃ separated analyzed the self-referential and other-referential trials. For the Self-referential
₁₃₇₄ trials, we found the main effect of valence, $F(1.97, 114.32) = 6.29$, $MSE = 367.25$,
₁₃₇₅ $p = .003$, $\hat{\eta}_G^2 = .006$; for the other-referential trials, the effect of valence is weaker,
₁₃₇₆ $F(1.95, 112.89) = 0.35$, $MSE = 699.50$, $p = .699$, $\hat{\eta}_G^2 = .001$. We then focused on the self
₁₃₇₇ conditions: the good-self condition (713 ± 12) is faster than neutral- (776 ± 11.8), $t(34) =$
₁₃₇₈ -7.396 , $p < .0001$, and bad-self (772 ± 10.1) conditions, $t(34) = -5.66$, $p < .0001$. But
₁₃₇₉ there is not difference between neutral- and bad-self conditions, $t(34) = 0.481$, $p = 0.881$.

₁₃₈₀ For the nonmatch trials, we didn't found any strong effect: identity, $F(1, 58) = 0.16$,
₁₃₈₁ $MSE = 1,547.37$, $p = .692$, $\hat{\eta}_G^2 = .000$, valence $F(1.96, 113.52) = 0.68$, $MSE = 390.26$,
₁₃₈₂ $p = .508$, $\hat{\eta}_G^2 = .000$, or interaction between the two $F(1.90, 110.27) = 0.04$,
₁₃₈₃ $MSE = 585.80$, $p = .953$, $\hat{\eta}_G^2 = .000$.

₁₃₈₄ **BGLM.**

₁₃₈₅ *Signal detection theory analysis of accuracy.*

₁₃₈₆ We found that the d prime is greater when shapes were associated with good self
₁₃₈₇ condition than with neutral self or bad self, but shapes associated with bad self and neutral

1388 self didn't show differences. comparing the self vs other under three condition revealed that
1389 shapes associated with good self is greater than with good other, but with a weak evidence.
1390 In contrast, for both neutral and bad valence condition, shapes associated with other had
1391 greater d prime than with self.

1392 *Reaction time.*

1393 In reaction times, we found that same trends in the match trials as in the RT: while
1394 the shapes associated with good self was greater than with good other (log mean diff =
1395 -0.02858, 95%HPD[-0.070898, 0.0154]), the direction is reversed for neutral and negative
1396 condition. see Figure 26

1397 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et
1398 al., 2020). We estimated separate drift rate (v), non-decision time (T_0), and boundary
1399 separation (a) for each condition. We found that the shapes tagged with good person has
1400 higher drift rate and higher boundary separation than shapes tagged with both neutral and
1401 bad person. Also, the shapes tagged with neutral person has a higher drift rate than
1402 shapes tagged with bad person, but not for the boundary separation. Finally, we found
1403 that shapes tagged with bad person had longer non-decision time (see figure 27)).

1404 **Experiment 4b: Morality as task-irrelevant variable**

1405 In study 4b, we changed the role of valence and identity in task. In this experiment,
1406 participants learn the association between moral valence and the made perceptual match
1407 judgments to associations between different moral valence and shapes as in study 1-3.
1408 Different from experiment 1 ~ 3, we made put the labels of "self/other" in the shapes so
1409 that identity served as an task irrelevant variable. As in experiment 4b, we also
1410 hypothesized that the intrinsic binding between morally good and self will enhance the
1411 performance of good self condition, even identity is irrelevant to the task.

1412 **Method.**

Participants.

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

Design.

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

Stimuli.

2 shapes were included (circle, square) and each appeared above a central fixation cross with the personal label appearing below. However, the shapes were not empty but with a two-Chinese-character word in the middle, the word was one of three labels with

1439 different moral valence: “good person”, “bad person” and “neutral person”. Before the
1440 experiment, participants learned the self/other association, and were informed to only
1441 response to the association between shapes’ configures and the labels below the fixation, but
1442 ignore the words within shapes. Besides the behavioral experiments, participants from
1443 Tsinghua community also finished questionnaires as Experiments 3, and participants from
1444 Wenzhou community finished a series of questionnaire as the other experiment finished in
1445 Wenzhou.

1446 ***Procedure.***

1447 The procedure was similar to Experiment 1. There were 6 blocks of trial, each with
1448 120 trials for 2017 data. Due to procedure error, the data collected in 2015 in Tsinghua
1449 community only have 60 trials for each block, i.e., 30 trials per condition.

1450 As in study 3a, before each task, the instruction showed the meaning of each label to
1451 participants. The self-matching task and other-matching task were randomized between
1452 participants. Each participant finished 6 blocks, each have 120 trials.

1453 ***Data Analysis.***

1454 Same as experiment 3a.

1455 ***Results.***

1456 ***NHST.***

1457 Figure 28 shows d prime and reaction times of experiment 3a. Less than 5% correct
1458 trials with less than 200ms reaction times were excluded.

1459 d prime.

1460 There was no evidence for the main effect of valence, $F(1.59, 69.94) = 2.34$,
1461 $MSE = 0.48$, $p = .115$, $\hat{\eta}_G^2 = .010$, but we found a main effect of self-relevance,
1462 $F(1, 44) = 0.00$, $MSE = 0.08$, $p = .994$, $\hat{\eta}_G^2 = .000$, as well as the interaction,
1463 $F(1.96, 86.41) = 0.53$, $MSE = 0.10$, $p = .585$, $\hat{\eta}_G^2 = .001$.

¹⁴⁶⁴ We then conducted separated ANOVA for self-referential and other-referential trials.

¹⁴⁶⁵ The valence effect was shown for the self-referential conditions, $F(1.75, 76.86) = 3.08$,

¹⁴⁶⁶ $MSE = 0.25$, $p = .058$, $\hat{\eta}_G^2 = .017$. Post-hoc test revealed that the Good-Self condition

¹⁴⁶⁷ (2.15 ± 0.12) was with greater d prime than Neutral condition $(1.83 \pm 0.12$, $t(34) = 3.36$,

¹⁴⁶⁸ $p = 0.0031$), and Bad-self condition (1.87 ± 0.12) , $t(34) = 2.955$, $p = 0.01$. There was

¹⁴⁶⁹ difference between neutral and bad condition, $t(34) = -0.039$, $p = 0.914$. However, no effect

¹⁴⁷⁰ of valence was found for the other-referential condition $F(1.63, 71.50) = 1.07$, $MSE = 0.33$,

¹⁴⁷¹ $p = .336$, $\hat{\eta}_G^2 = .006$.

¹⁴⁷² *Reaction time.*

¹⁴⁷³ We found interaction between Matchness and Valence ($F(1.87, 82.50) = 18.58$,

¹⁴⁷⁴ $MSE = 1,291.12$, $p < .001$, $\hat{\eta}_G^2 = .023$) and then analyzed the matched trials and

¹⁴⁷⁵ nonmatch trials separately, as in previous experiments.

¹⁴⁷⁶ For the match trials, we found that the interaction between identity and valence,

¹⁴⁷⁷ $F(1.86, 81.84) = 5.22$, $MSE = 308.30$, $p = .009$, $\hat{\eta}_G^2 = .003$, as well as the main effect of

¹⁴⁷⁸ valence $F(1.80, 79.37) = 11.04$, $MSE = 2,937.54$, $p < .001$, $\hat{\eta}_G^2 = .059$, but not the effect of

¹⁴⁷⁹ identity $F(1, 44) = 0.23$, $MSE = 263.26$, $p = .632$, $\hat{\eta}_G^2 = .000$. As for the d prime, we

¹⁴⁸⁰ separated analyzed the self-referential and other-referential trials. For the Self-referential

¹⁴⁸¹ trials, we found the main effect of valence, $F(1.74, 76.48) = 13.69$, $MSE = 1,732.08$,

¹⁴⁸² $p < .001$, $\hat{\eta}_G^2 = .079$; for the other-referential trials, the effect of valence is weaker,

¹⁴⁸³ $F(1.87, 82.44) = 7.09$, $MSE = 1,527.43$, $p = .002$, $\hat{\eta}_G^2 = .043$. We then focused on the self

¹⁴⁸⁴ conditions: the good-self condition (713 ± 12) is faster than neutral- (776 ± 11.8) , $t(34) =$

¹⁴⁸⁵ -7.396 , $p < .0001$, and bad-self (772 ± 10.1) conditions, $t(34) = -5.66$, $p < .0001$. But

¹⁴⁸⁶ there is not difference between neutral- and bad-self conditions, $t(34) = 0.481$, $p = 0.881$.

¹⁴⁸⁷ For the nonmatch trials, we didn't found any strong effect: identity, $F(1, 44) = 1.96$,

¹⁴⁸⁸ $MSE = 319.47$, $p = .169$, $\hat{\eta}_G^2 = .001$, valence $F(1.69, 74.54) = 6.59$, $MSE = 886.19$,

¹⁴⁸⁹ $p = .004$, $\hat{\eta}_G^2 = .010$, or interaction between the two $F(1.88, 82.57) = 0.31$, $MSE = 316.96$,

₁₄₉₀ $p = .718$, $\hat{\eta}_G^2 = .000$.

₁₄₉₁ **BGLM.**

₁₄₉₂ *Signal detection theory analysis of accuracy.*

₁₄₉₃ We found that the d prime is greater when shapes were associated with good self
₁₄₉₄ condition than with neutral self or bad self, but shapes associated with bad self and neutral
₁₄₉₅ self didn't show differences. comparing the self vs other under three condition revealed that
₁₄₉₆ shapes associated with good self is greater than with good other, but with a weak evidence.
₁₄₉₇ In contrast, for both neutral and bad valence condition, shapes associated with other had
₁₄₉₈ greater d prime than with self.

₁₄₉₉ *Reaction time.*

₁₅₀₀ In reaction times, we found that same trends in the match trials as in the RT: while
₁₅₀₁ the shapes associated with good self was greater than with good other ($\log \text{mean diff} =$
₁₅₀₂ -0.02858 , $95\% \text{HPD}[-0.070898, 0.0154]$), the direction is reversed for neutral and negative
₁₅₀₃ condition. see Figure 29

₁₅₀₄ **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et
₁₅₀₅ al., 2020). We estimated separate drift rate (v), non-decision time (T_0), and boundary
₁₅₀₆ separation (a) for each condition. We found that the shapes tagged with good person has
₁₅₀₇ higher drift rate and higher boundary separation than shapes tagged with both neutral and
₁₅₀₈ bad person. Also, the shapes tagged with neutral person has a higher drift rate than
₁₅₀₉ shapes tagged with bad person, but not for the boundary separation. Finally, we found
₁₅₁₀ that shapes tagged with bad person had longer non-decision time (see figure 30)).

1511

Results

1512 **Effect of moral valence**

1513 In this part, we synthesized results from experiment 1a, 1b, 1c, 2, 5 and 6a. Data
1514 from 192 participants were included in these analyses. We found differences between
1515 positive and negative conditions on RT was Cohen's $d = -0.58 \pm 0.06$, 95% CI [-0.70 -0.47];
1516 on d' was Cohen's $d = 0.24 \pm 0.05$, 95% CI [0.15 0.34]. The effect was also observed
1517 between positive and neutral condition, RT: Cohen's $d = -0.44 \pm 0.10$, 95% CI [-0.63
1518 -0.25]; d' : Cohen's $d = 0.31 \pm 0.07$, 95% CI [0.16 0.45]. And the difference between neutral
1519 and bad conditions are not significant, RT: Cohen's $d = 0.15 \pm 0.07$, 95% CI [0.00 0.30];
1520 d' : Cohen's $d = 0.07 \pm 0.07$, 95% CI [-0.08 0.21]. See Figure 31 left panel.

1521 **Interaction between valence and self-reference**

1522 In this part, we combined the experiments that explicitly manipulated the
1523 self-reference and valence, which includes 3a, 3b, 6b, 7a, and 7b. For the positive versus
1524 negative contrast, data were from five experiments with 178 participants; for positive
1525 versus neutral and neutral versus negative contrasts, data were from three experiments ((

1526 3a, 3b, and 6b) with 108 participants.

1527 In most of these experiments, the interaction between self-reference and valence was
1528 significant (see results of each experiment in supplementary materials). In the
1529 mini-meta-analysis, we analyzed the valence effect for self-referential condition and
1530 other-referential condition separately.

1531 For the self-referential condition, we found the same pattern as in the first part of
1532 results. That is we found significant differences between positive and neutral as well as
1533 positive and negative, but not neutral and negative. The effect size of RT between positive
1534 and negative is Cohen's $d = -0.89 \pm 0.12$, 95% CI [-1.11 -0.66]; on d' was Cohen's $d = 0.61$

1535 ± 0.09 , 95% CI [0.44 0.78]. The effect was also observed between positive and neutral
1536 condition, RT: Cohen's $d = -0.76 \pm 0.13$, 95% CI [-1.01 -0.50]; d' : Cohen's $d = 0.69 \pm$
1537 0.14, 95% CI [0.42 0.96]. And the difference between neutral and bad conditions are not
1538 significant, RT: Cohen's $d = 0.03 \pm 0.13$, 95% CI [-0.22 0.29]; d' : Cohen's $d = 0.08 \pm 0.08$,
1539 95% CI [-0.07 0.24]. See Figure 31 the middle panel.

1540 For the other-referential condition, we found that only the difference between positive
1541 and negative on RT was significant, all the other conditions were not. The effect size of RT
1542 between positive and negative is Cohen's $d = -0.28 \pm 0.05$, 95% CI [-0.38 -0.17]; on d' was
1543 Cohen's $d = -0.02 \pm 0.08$, 95% CI [-0.17 0.13]. The effect was not observed between
1544 positive and neutral condition, RT: Cohen's $d = -0.12 \pm 0.10$, 95% CI [-0.31 0.06]; d' :
1545 Cohen's $d = 0.01 \pm 0.08$, 95% CI [-0.16 0.17]. And the difference between neutral and bad
1546 conditions are not significant, RT: Cohen's $d = 0.14 \pm 0.09$, 95% CI [-0.03 0.31]; d' :
1547 Cohen's $d = 0.05 \pm 0.07$, 95% CI [-0.08 0.18]. See Figure 31 right panel.

1548 Generalizability of the valence effect

1549 In this part, we reported the results from experiment 4 in which either moral valence
1550 or self-reference were manipulated as task-irrelevant stimuli.

1551 For experiment 4a, when self-reference was the target and moral valence was
1552 task-irrelevant, we found that only under the implicit self-referential condition, i.e., when
1553 the moral words were presented as task irrelevant stimuli, there was the main effect of
1554 valence and interaction between valence and reference for both d prime and RT (See
1555 supplementary results for the detailed statistics). For d prime, we found good-self
1556 condition (2.55 ± 0.86) had higher d prime than bad-self condition (2.38 ± 0.80); good self
1557 condition was also higher than neutral self (2.45 ± 0.78) but there was not statistically
1558 significant, while the neutral-self condition was higher than bad self condition and not
1559 significant neither. For reaction times, good-self condition (654.26 ± 67.09) were faster

1560 relative to bad-self condition (665.64 ± 64.59), and over neutral-self condition ($664.26 \pm$
1561 64.71). The difference between neutral-self and bad-self conditions were not significant.
1562 However, for the other-referential condition, there was no significant differences between
1563 different valence conditions. See Figure 32.

1564 For experiment 4b, when valence was the target and the identity was task-irrelevant,
1565 we found a strong valence effect (see supplementary results and Figure 33, Figure 34).

1566 In this experiment, the advantage of good-self condition can only be disentangled by
1567 comparing the self-referential and other-referential conditions. Therefore, we calculated the
1568 differences between the valence effect under self-referential and other referential conditions
1569 and used the weighted variance as the variance of this differences. We found this
1570 modulation effect on RT. The valence effect of RT was stronger in self-referential than
1571 other-referential for the Good vs. Neutral condition (-0.33 ± 0.01), and to a less extent the
1572 Good vs. Bad condition (-0.17 ± 0.01). While the size of the other effect's CI included
1573 zero, suggestion those effects didn't differ from zero. See Figure 35.

1574 Specificity of valence effect

1575 In this part, we analyzed the results from experiment 5, which included positive,
1576 neutral, and negative valence from four different domains: morality, emotion, aesthetics of
1577 human, and aesthetics of scene. We found interaction between valence and domain for both
1578 *d* prime and RT (match trials). A common pattern appeared in all four domains: each
1579 domain showed a binary results instead of gradient on both *d* prime and RT. For morality,
1580 aesthetics of human, and aesthetics of scene, there was a positivity effect where the positive
1581 conditions had advantages over both neutral (greater *d* prime and faster RT), while neutral
1582 and negative conditions didn't differ from each other. But for the emotional stimuli, there
1583 was a reversed negativity effect: positive and neutral conditions were not significantly
1584 different from each other but both had advantage over negative conditions. See

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

Self-reported personal distance

See Figure 37.

Correlation analyses

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

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

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

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

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

1609 ($r = 0.1571$, 95% CI [], $p = 0.031$). At the same time, the non-decision time is negatively
1610 correlated with Self-Bad distance ($r = 0.169$, 95% CI [], $p = 0.0197$). See Figure 38.

1611 Second, we found the boundary separation of neutral condition is positively
1612 correlated with the personal distance between self and good distance ($r = 0.189$, 95% CI [],
1613 $p = 0.036$), but negatively correlated with self-neutral distance($r = -0.183$, 95% CI [],
1614 $p = 0.042$). Also, the drift rate of the neutral condition is positively correlated with the
1615 Self-Bad distance ($r = 0.177$, 95% CI [], $p = 0.048$).a. See figure 39

1616 We also explored the correlation between behavioral data and questionnaire scores
1617 separately for experiments with and without self-referential, however, the sample size is
1618 very low for some conditions.

1619 **Discussion**

1620 **References**

- 1621 Abele, A. E., Ellemers, N., Fiske, S. T., Koch, A., & Yzerbyt, V. (2020). Navigating the
1622 social world: Toward an integrated framework for evaluating self, individuals, and
1623 groups. *Psychological Review*. <https://doi.org/10.1037/rev0000262>
- 1624 Amodio, D. M. (2019). Social cognition 2.0: An interactive memory systems account.
1625 *Trends in Cognitive Sciences*, 23(1), 21–33.
1626 <https://doi.org/10.1016/j.tics.2018.10.002>
- 1627 Anderson, E., Siegel, E. H., Bliss-Moreau, E., & Barrett, L. F. (2011). The visual impact
1628 of gossip. *Science*, 332(6036), 1446–1448. <https://doi.org/10.1126/science.1201574>
- 1629 Brainard, D. H. (1997). The psychophysics toolbox. *Spatial Vision*, 10(4), 433–436.
1630 Journal Article.
- 1631 Bürkner, P.-C. (2017). Brms: An r package for bayesian multilevel models using stan.
1632 *Journal of Statistical Software*; Vol 1, Issue 1 (2017). Journal Article. Retrieved

- 1633 from
1634 <https://www.jstatsoft.org/v080/i01> <http://dx.doi.org/10.18637/jss.v080.i01>
- 1635 Carlson, R. W., Maréchal, M. A., Oud, B., Fehr, E., & Crockett, M. J. (2020). Motivated
1636 misremembering of selfish decisions. *Nature Communications*, 11(1), 2100.
1637 <https://doi.org/10.1038/s41467-020-15602-4>
- 1638 Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M., ...
1639 Riddell, A. (2017). Stan: A probabilistic programming language. *Journal of
1640 Statistical Software*, 76(1). Journal Article. <https://doi.org/10.18637/jss.v076.i01>
- 1641 Cooper, H., Hedges, L. V., & Valentine, J. C. (2009). *The handbook of research synthesis
1642 and meta-analysis* (2nd ed.). Book, New York: Sage.
- 1643 DeCarlo, L. T. (1998). Signal detection theory and generalized linear models. *Psychological
1644 Methods*, 3(2), 186–205. Journal Article. <https://doi.org/10.1037/1082-989X.3.2.186>
- 1645 Dzhelyova, M., Perrett, D. I., & Jentzsch, I. (2012). Temporal dynamics of trustworthiness
1646 perception. *Brain Research*, 1435, 81–90.
1647 <https://doi.org/10.1016/j.brainres.2011.11.043>
- 1648 Ellemers, N., Toorn, J. van der, Paunov, Y., & Leeuwen, T. van. (2019). The psychology of
1649 morality: A review and analysis of empirical studies published from 1940 through
1650 2017. *Personality and Social Psychology Review*, 23(4), 332–366.
1651 <https://doi.org/10.1177/1088868318811759>
- 1652 Enock, F. E., Hewstone, M. R. C., Lockwood, P. L., & Sui, J. (2020). Overlap in
1653 processing advantages for minimal ingroups and the self. *Scientific Reports*, 10(1),
1654 18933. <https://doi.org/10.1038/s41598-020-76001-9>
- 1655 Enock, F., Sui, J., Hewstone, M., & Humphreys, G. W. (2018). Self and team prioritisation
1656 effects in perceptual matching: Evidence for a shared representation. *Acta
1657 Psychologica*, 182, 107–118. <https://doi.org/10.1016/j.actpsy.2017.11.011>

- 1658 Faul, F., Erdfelder, E., Buchner, A., & Lang, A.-G. (2009). Statistical power analyses using
1659 g*Power 3.1: Tests for correlation and regression analyses. *Behavior Research
1660 Methods*, 41(4), 1149–1160. <https://doi.org/10.3758/BRM.41.4.1149>
- 1661 Firestone, C., & Scholl, B. J. (2015). Enhanced visual awareness for morality and pajamas?
1662 Perception vs. Memory in “top-down” effects. *Cognition*, 136, 409–416.
1663 <https://doi.org/10.1016/j.cognition.2014.10.014>
- 1664 Freeman, J. B., & Ambady, N. (2011). A dynamic interactive theory of person construal.
1665 *Psychological Review*, 118(2), 247–279. <https://doi.org/10.1037/a0022327>
- 1666 Gantman, A. P., & Van Bavel, J. J. (2014). The moral pop-out effect: Enhanced
1667 perceptual awareness of morally relevant stimuli. *Cognition*, 132(1), 22–29.
1668 <https://doi.org/10.1016/j.cognition.2014.02.007>
- 1669 Goh, J. X., Hall, J. A., & Rosenthal, R. (2016). Mini meta-analysis of your own studies:
1670 Some arguments on why and a primer on how. *Social and Personality Psychology
1671 Compass*, 10(10), 535–549. Journal Article. <https://doi.org/10.1111/spc3.12267>
- 1672 Goodwin, G. P. (2015). Moral character in person perception. *Current Directions in
1673 Psychological Science*, 24(1), 38–44. <https://doi.org/10.1177/0963721414550709>
- 1674 Goodwin, G. P., Piazza, J., & Rozin, P. (2014). Moral character predominates in person
1675 perception and evaluation. *Journal of Personality and Social Psychology*, 106(1),
1676 148–168. <https://doi.org/10.1037/a0034726>
- 1677 Henrich, J., Heine, S. J., & Norenzayan, A. (2010). The weirdest people in the world?
1678 *Behavioral and Brain Sciences*, 33(2), 61–83.
1679 <https://doi.org/10.1017/S0140525X0999152X>
- 1680 Hofmann, W., Wisneski, D. C., Brandt, M. J., & Skitka, L. J. (2014). Morality in everyday
1681 life. *Science*, 345(6202), 1340–1343. <https://doi.org/10.1126/science.1251560>
- 1682 Hu, C.-P., Lan, Y., Macrae, C. N., & Sui, J. (2020). Good me bad me: Does valence

- 1683 influence self-prioritization during perceptual decision-making? *Collabra: Psychology*, 6(1), 20. Journal Article. <https://doi.org/10.1525/collabra.301>
- 1684
- 1685 Juechems, K., & Summerfield, C. (2019). Where does value come from? *Trends in Cognitive Sciences*, 23(10), 836–850. <https://doi.org/10.1016/j.tics.2019.07.012>
- 1686
- 1687 Kirby, K. N., & Gerlanc, D. (2013). BootES: An r package for bootstrap confidence intervals on effect sizes. *Behavior Research Methods*, 45(4), 905–927. <https://doi.org/10.3758/s13428-013-0330-5>
- 1688
- 1689
- 1690 Klein, N., & Epley, N. (2016). Maybe holier, but definitely less evil, than you: Bounded self-righteousness in social judgment. *Journal of Personality and Social Psychology*, 110(5), 660–674. <https://doi.org/10.1037/pspa0000050>
- 1691
- 1692
- 1693 Liu, Q., Wang, F., Yan, W., Peng, K., Sui, J., & Hu, C.-P. (2020). Questionnaire data from the revision of a chinese version of free will and determinism plus scale. *Journal of Open Psychology Data*, 8(1), 1. Journal Article. <https://doi.org/10.5334/jopd.49/>
- 1694
- 1695
- 1696 Matzke, D., & Wagenmakers, E.-J. (2009). Psychological interpretation of the ex-gaussian and shifted wald parameters: A diffusion model analysis. *Psychonomic Bulletin & Review*, 16(5), 798–817. <https://doi.org/10.3758/PBR.16.5.798>
- 1697
- 1698
- 1699 McHugh, C., McGann, M., Igou, E. R., & Kinsella, E. (2019). *Moral judgment as categorization (MJAC)*. PsyArXiv. <https://doi.org/10.31234/osf.io/72dzp>
- 1700
- 1701 McNeish, D., & Wolf, M. G. (2020). Thinking twice about sum scores. *Behavior Research Methods*. <https://doi.org/10.3758/s13428-020-01398-0>
- 1702
- 1703 Monin, B., & Jordan, A. H. (2009). The dynamic moral self: A social psychological perspective. In *Personality, identity, and character: Explorations in moral psychology* (pp. 341–354). New York, NY, US: Cambridge University Press.
- 1704
- 1705
- 1706 <https://doi.org/10.1017/CBO9780511627125.016>
- 1707 Pelli, D. G. (1997). The videotoolbox software for visual psychophysics: Transforming

- 1708 numbers into movies. *Spatial Vision*, 10(4), 437–442. Journal Article.
- 1709 Reicher, S., & Hopkins, N. (2016). Perception, action, and the social dynamics of the
1710 variable self. *Psychological Inquiry*, 27(4), 341–347.
1711 <https://doi.org/10.1080/1047840X.2016.1217584>
- 1712 Rouder, J. N., & Lu, J. (2005). An introduction to bayesian hierarchical models with an
1713 application in the theory of signal detection. *Psychonomic Bulletin & Review*,
1714 12(4), 573–604. Journal Article. <https://doi.org/10.3758/bf03196750>
- 1715 Rousselet, G. A., & Wilcox, R. R. (2019). Reaction times and other skewed distributions:
1716 Problems with the mean and the median. *Meta-Psychology*. preprint.
1717 <https://doi.org/10.1101/383935>
- 1718 Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2013). Life after p-hacking. Conference
1719 Proceedings. <https://doi.org/10.2139/ssrn.2205186>
- 1720 Stanislaw, H., & Todorov, N. (1999). Calculation of signal detection theory measures.
1721 *Behavior Research Methods, Instruments, & Computers*, 31(1), 137–149. Journal
1722 Article. <https://doi.org/10.3758/BF03207704>
- 1723 Stanley, M. L., & De Brigard, F. (2019). Moral memories and the belief in the good self.
1724 *Current Directions in Psychological Science*, 28(4), 387–391.
1725 <https://doi.org/10.1177/0963721419847990>
- 1726 Stein, T., Grubb, C., Bertrand, M., Suh, S. M., & Verosky, S. C. (2017). No impact of
1727 affective person knowledge on visual awareness: Evidence from binocular rivalry and
1728 continuous flash suppression. *Emotion*, 17(8), 1199–1207.
1729 <https://doi.org/10.1037/emo0000305>
- 1730 Stolier, R. M., & Freeman, J. B. (2016). Functional and temporal considerations for
1731 top-down influences in social perception. *Psychological Inquiry*, 27(4), 352–357.
1732 <https://doi.org/10.1080/1047840X.2016.1216034>

- 1733 Strohminger, N., Knobe, J., & Newman, G. (2017). The true self: A psychological concept
1734 distinct from the self: *Perspectives on Psychological Science*.
1735 <https://doi.org/10.1177/1745691616689495>
- 1736 Sui, J., He, X., & Humphreys, G. W. (2012). Perceptual effects of social salience: Evidence
1737 from self-prioritization effects on perceptual matching. *Journal of Experimental
1738 Psychology: Human Perception and Performance*, 38(5), 1105–1117. Journal
1739 Article. <https://doi.org/10.1037/a0029792>
- 1740 Tappin, B. M., & McKay, R. T. (2017). The illusion of moral superiority. *Social
1741 Psychological and Personality Science*, 8(6), 623–631.
1742 <https://doi.org/10.1177/1948550616673878>
- 1743 Turner, J. C., Hogg, M. A., Oakes, P. J., Reicher, S. D., & Wetherell, M. S. (1987).
1744 *Rediscovering the social group: A self-categorization theory*. Cambridge, MA, US:
1745 Basil Blackwell.
- 1746 Turner, J. C., Oakes, P. J., Haslam, S. A., & McGarty, C. (1994). Self and collective:
1747 Cognition and social context. *Personality and Social Psychology Bulletin*, 20(5),
1748 454–463. <https://doi.org/10.1177/0146167294205002>
- 1749 Uhlmann, E. L., Pizarro, D. A., & Diermeier, D. (2015). A person-centered approach to
1750 moral judgment: *Perspectives on Psychological Science*.
1751 <https://doi.org/10.1177/1745691614556679>
- 1752 Verosky, S. C., & Todorov, A. (2010). Differential neural responses to faces physically
1753 similar to the self as a function of their valence. *NeuroImage*, 49(2), 1690–1698.
1754 <https://doi.org/10.1016/j.neuroimage.2009.10.017>
- 1755 Waytz, A., Dungan, J., & Young, L. (2013). The whistleblower's dilemma and the
1756 fairness–loyalty tradeoff. *Journal of Experimental Social Psychology*, 49(6),
1757 1027–1033. <https://doi.org/10.1016/j.jesp.2013.07.002>

- 1758 Wiecki, T. V., Sofer, I., & Frank, M. J. (2013). HDDM: Hierarchical bayesian estimation of
1759 the drift-diffusion model in python. *Frontiers in Neuroinformatics*, 7.
1760 <https://doi.org/10.3389/fninf.2013.00014>
- 1761 Willis, J., & Todorov, A. (2006). First impressions: Making up your mind after a 100-ms
1762 exposure to a face. *Psychological Science*, 17(7), 592–598.
1763 <https://doi.org/10.1111/j.1467-9280.2006.01750.x>
- 1764 Xiao, Y. J., Coppin, G., & Bavel, J. J. V. (2016). Perceiving the world through
1765 group-colored glasses: A perceptual model of intergroup relations. *Psychological
1766 Inquiry*, 27(4), 255–274. <https://doi.org/10.1080/1047840X.2016.1199221>

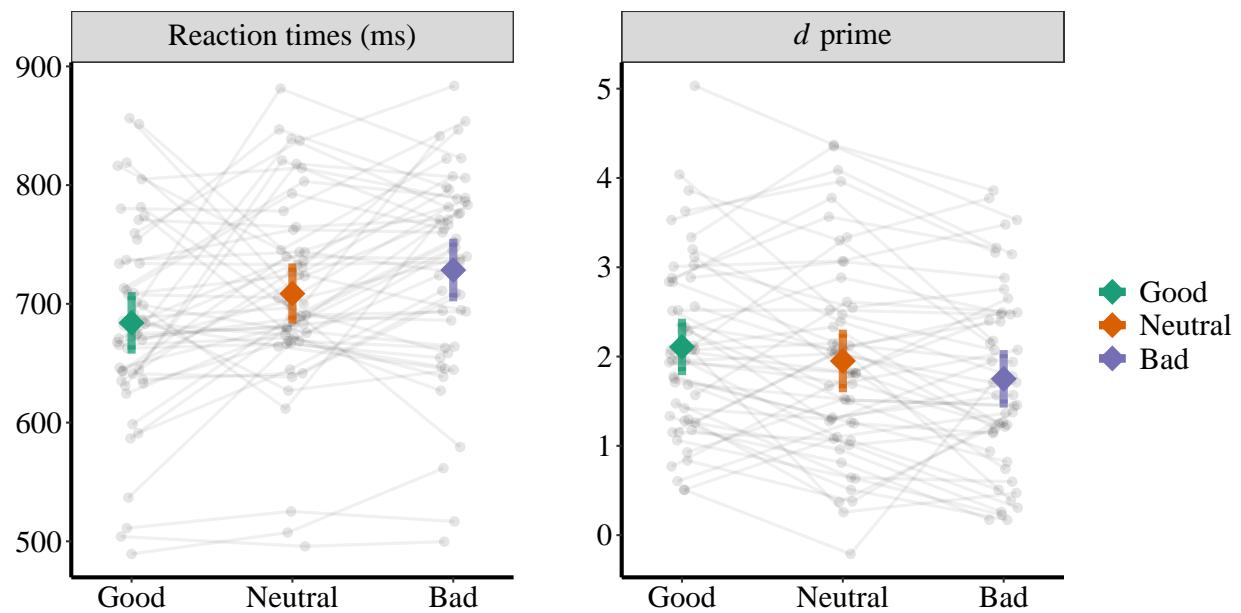


Figure 1. RT and d' prime of Experiment 1a.

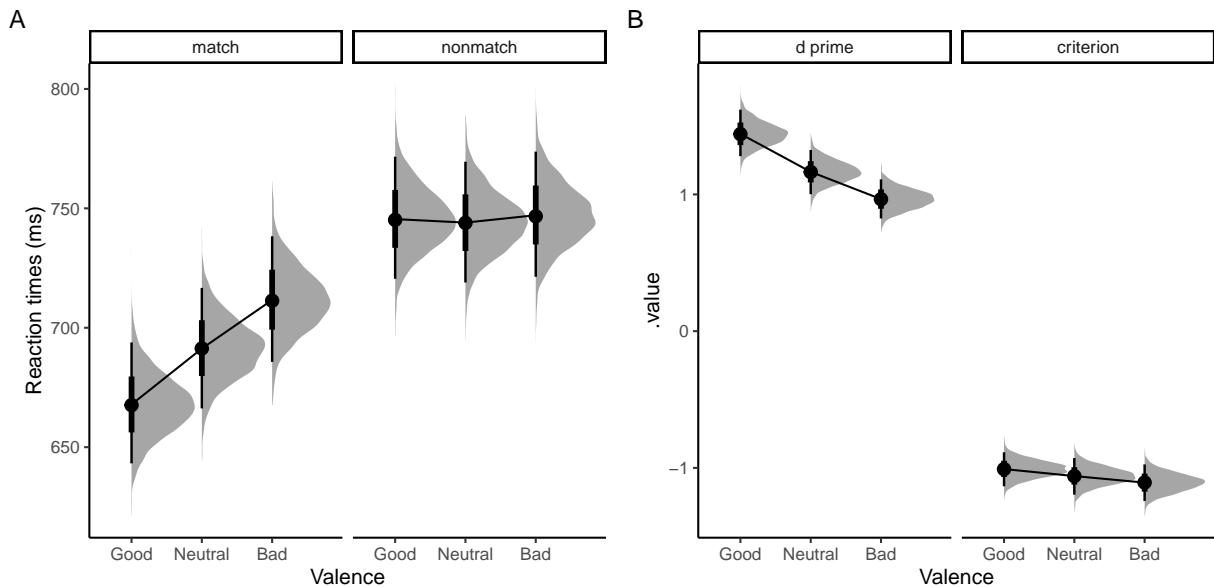


Figure 2. Exp1a: Results of Bayesian GLM analysis.

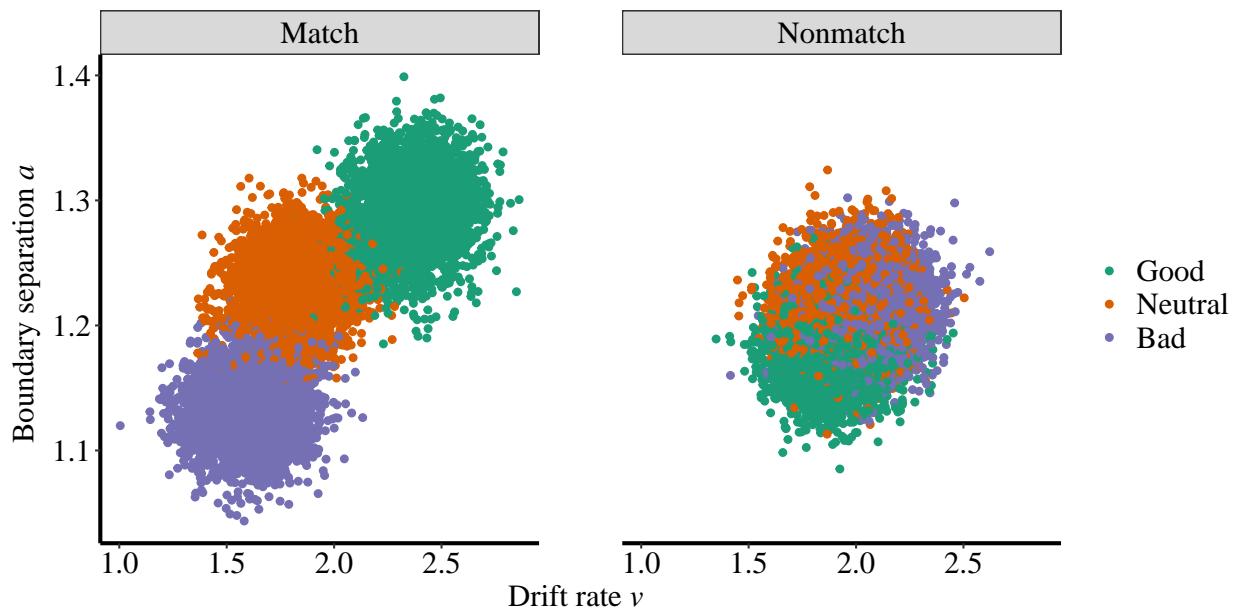


Figure 3. Exp1a: Results of HDDM.

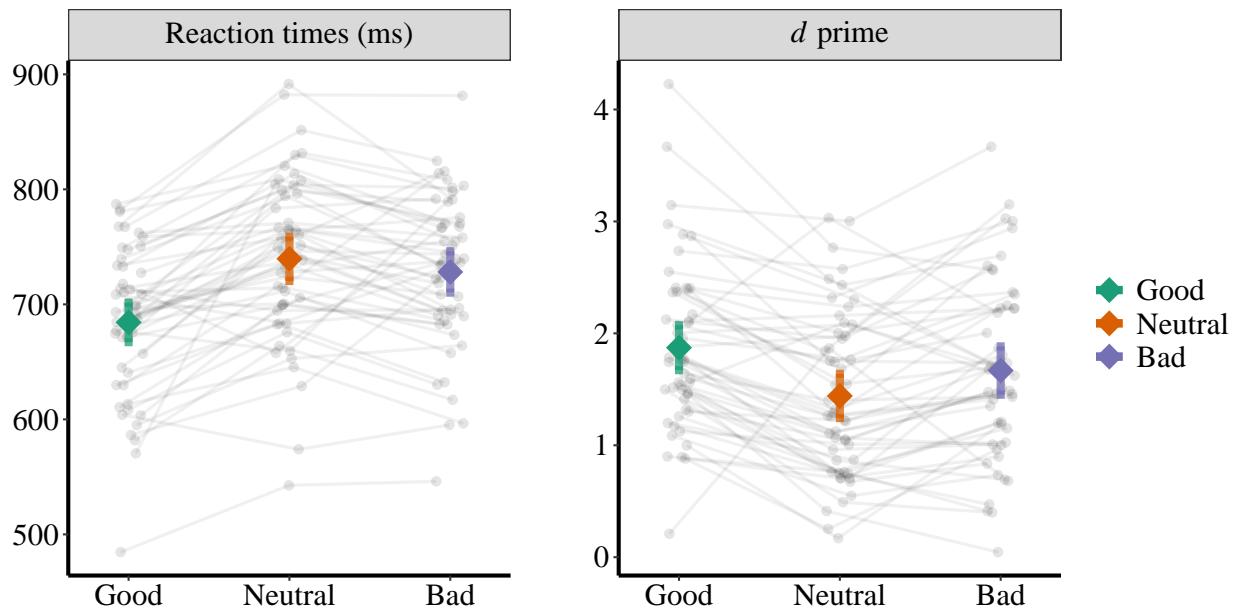


Figure 4. RT and d' of Experiment 1b.

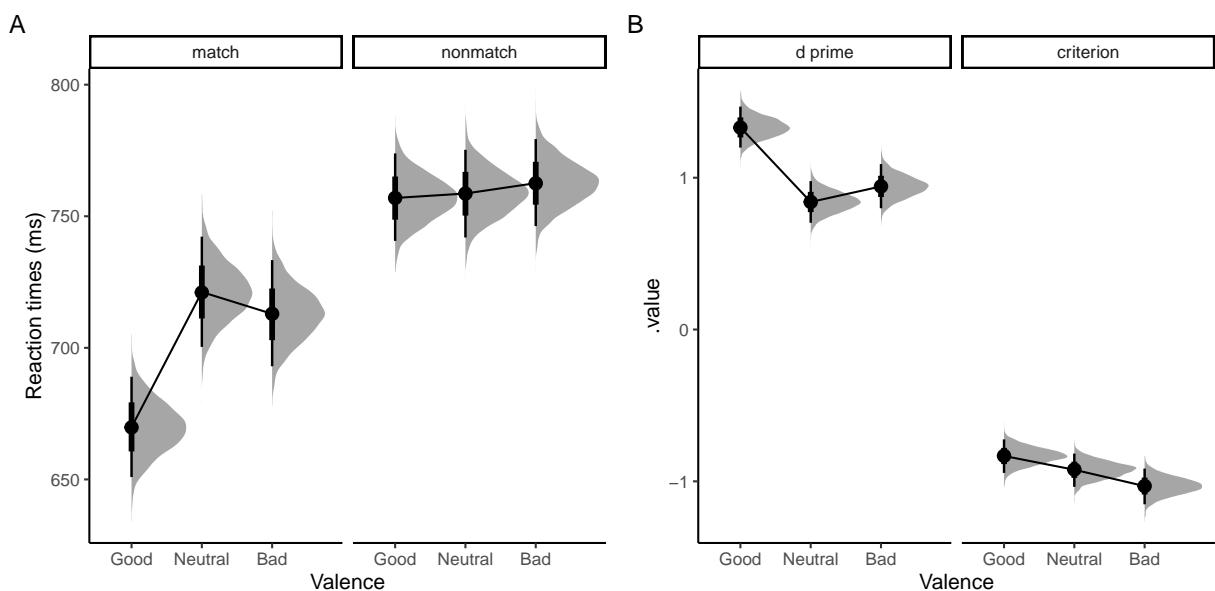


Figure 5. Exp1b: Results of Bayesian GLM analysis.

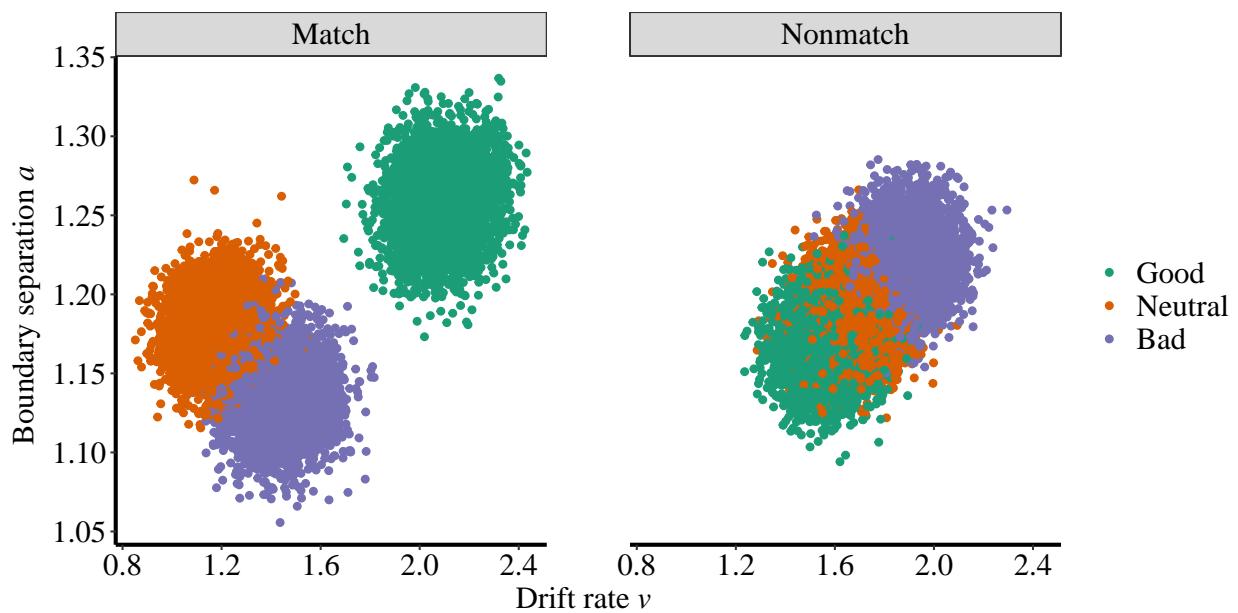


Figure 6. Exp1b: Results of HDDM.

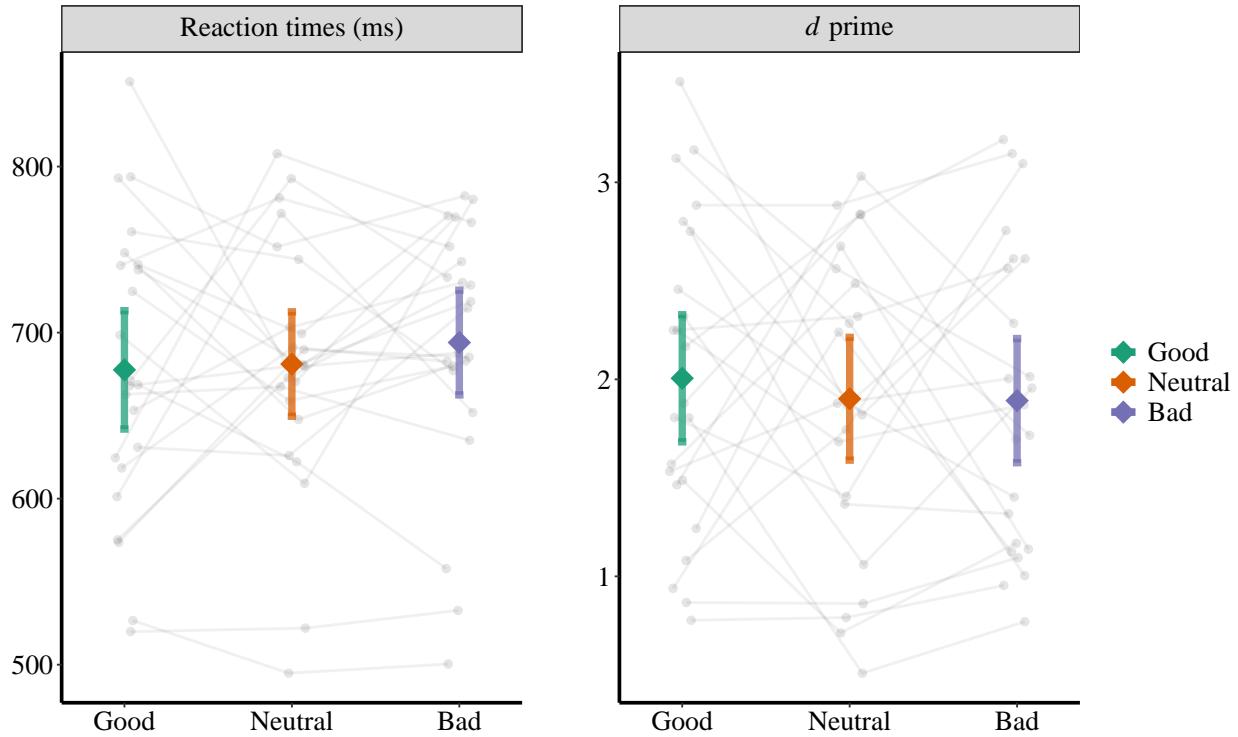


Figure 7. RT and d' prime of Experiment 1c.

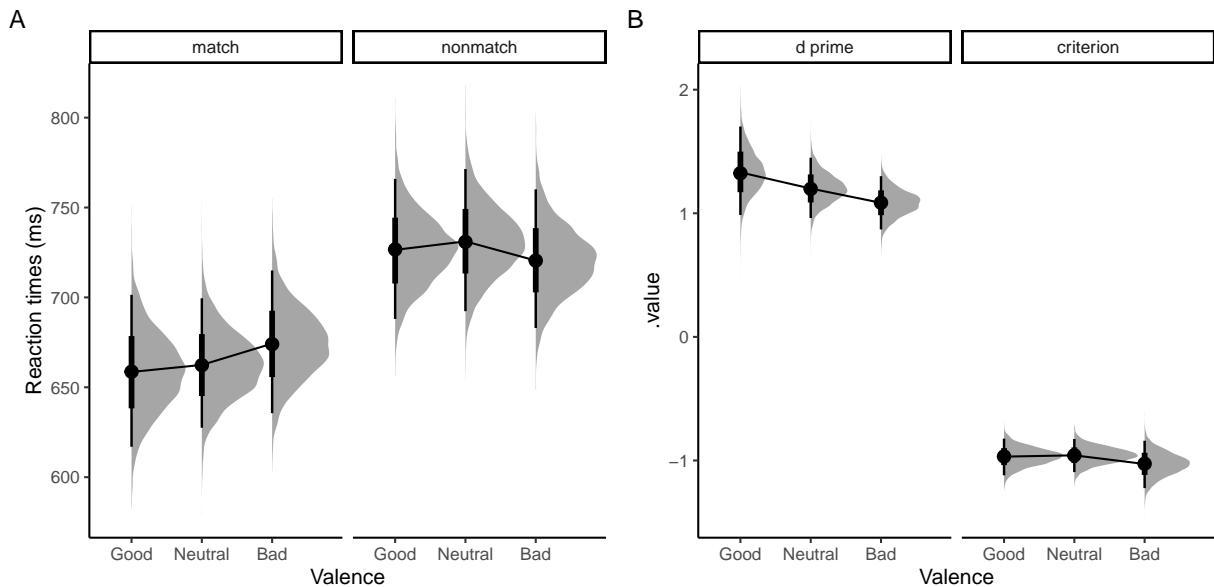


Figure 8. Exp1c: Results of Bayesian GLM analysis.

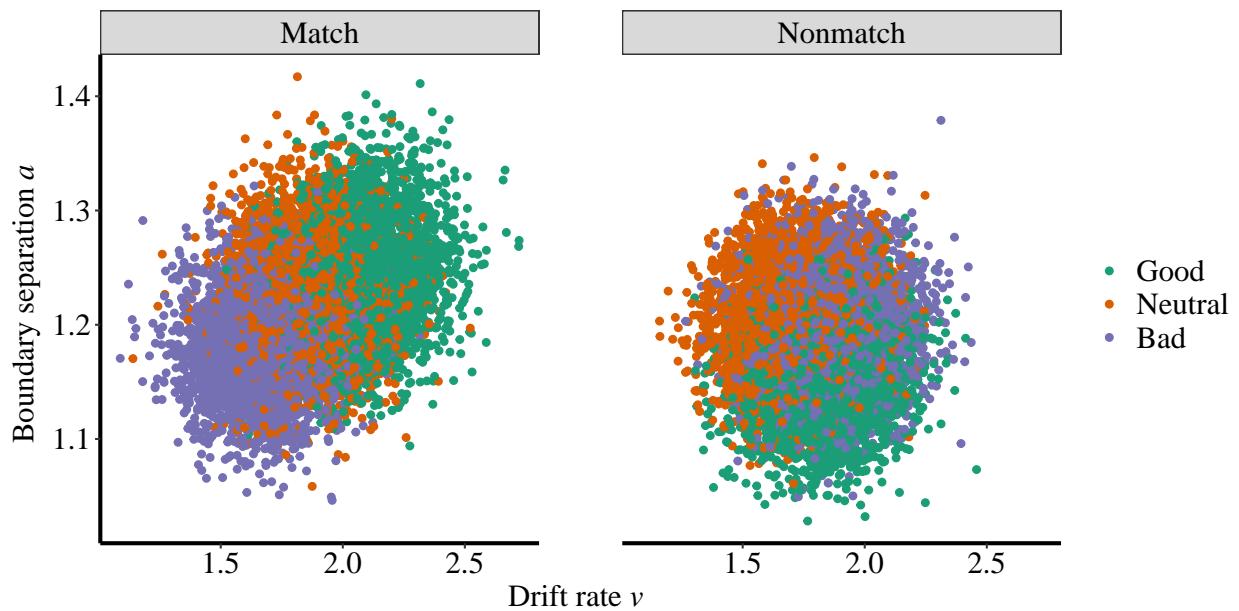


Figure 9. Exp1c: Results of HDDM.

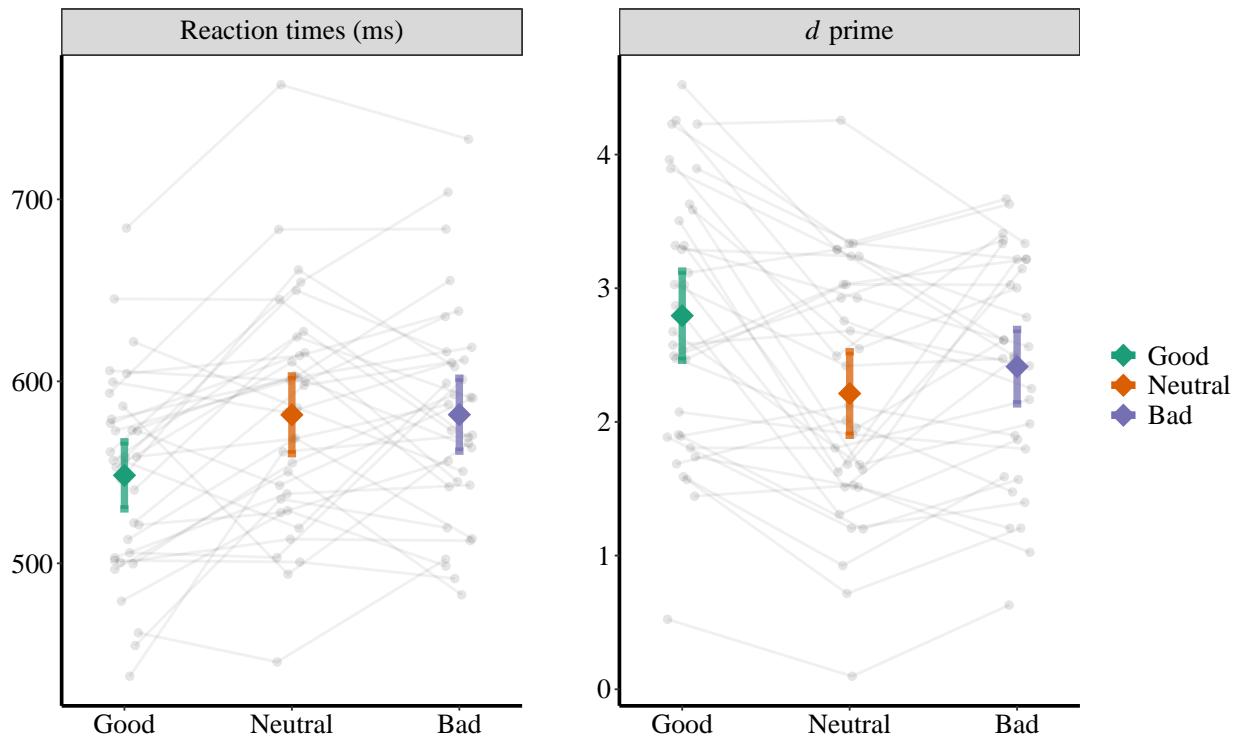


Figure 10. RT and d' of Experiment 2.

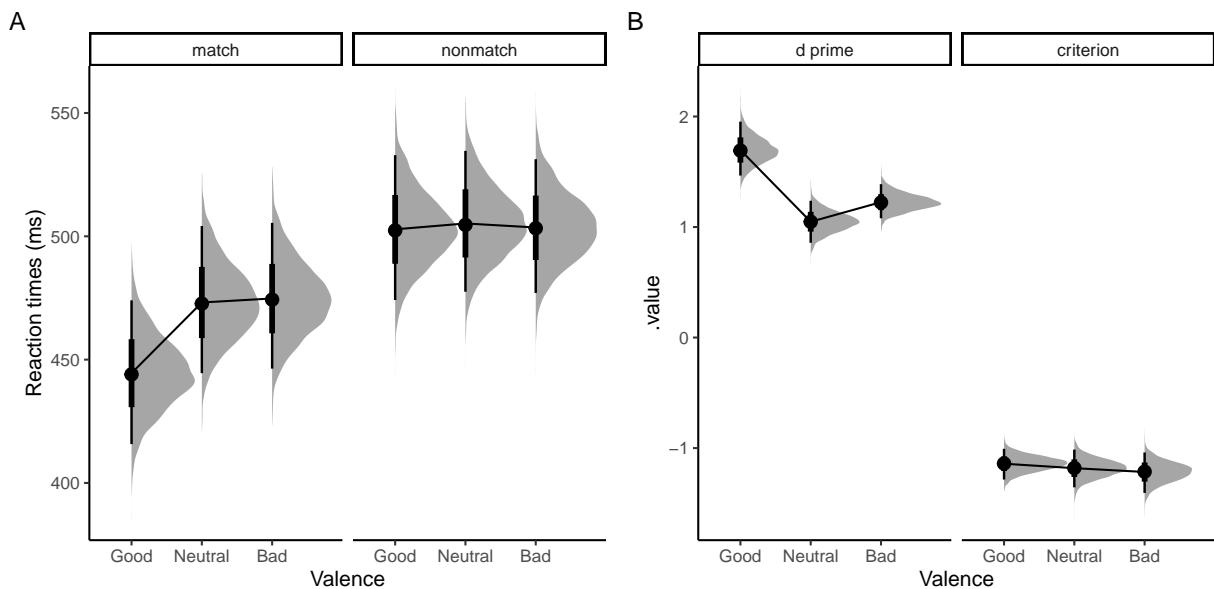


Figure 11. Exp2: Results of Bayesian GLM analysis.

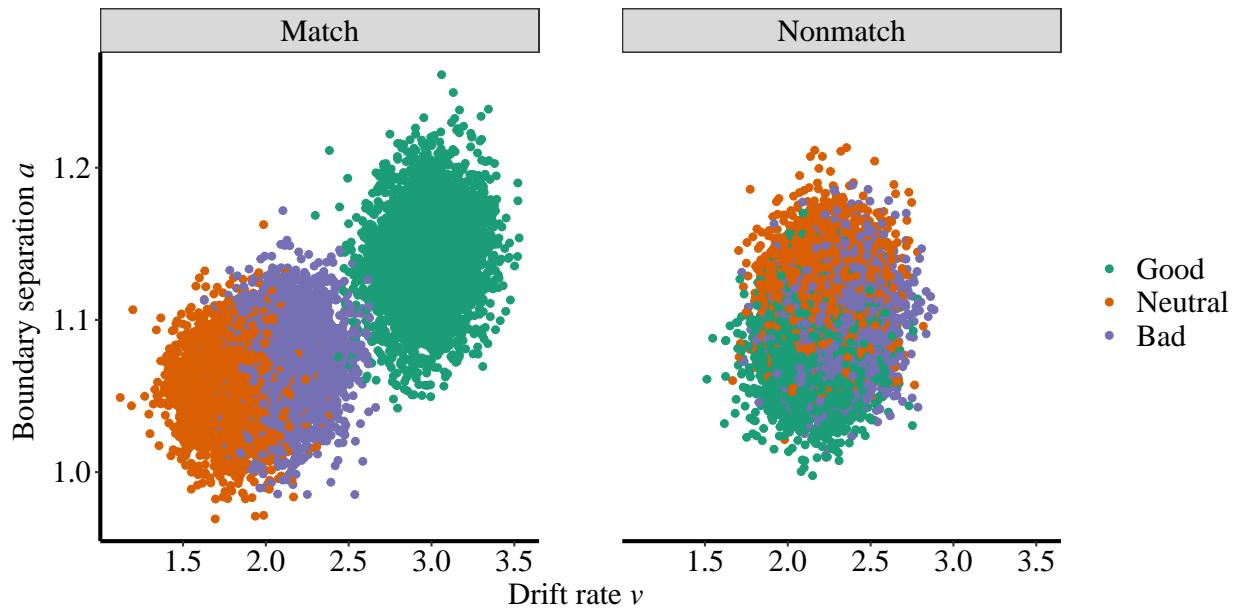


Figure 12. Exp2: Results of HDDM.

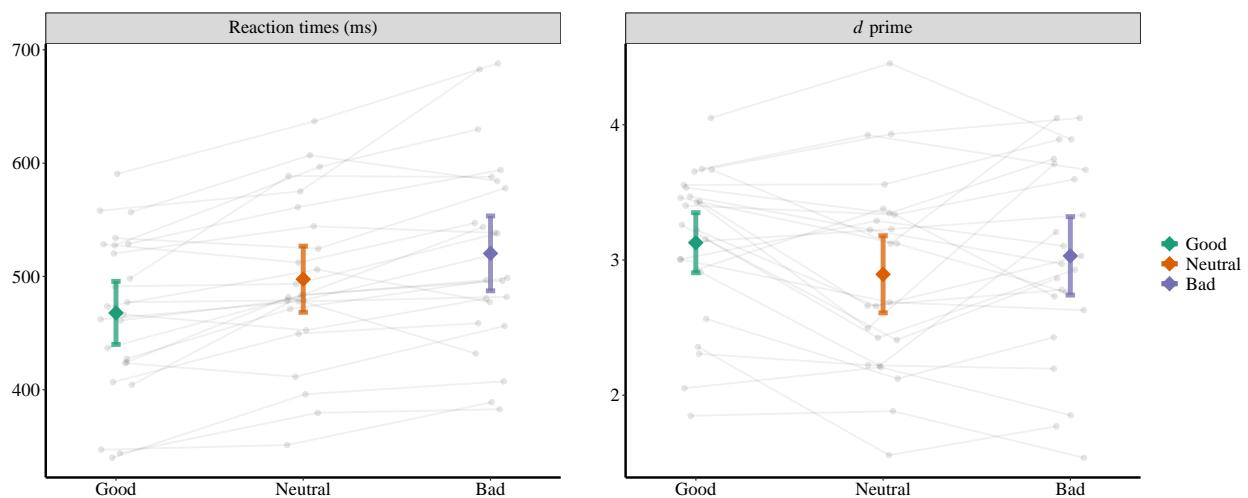


Figure 13. RT and d' of Experiment 6a.

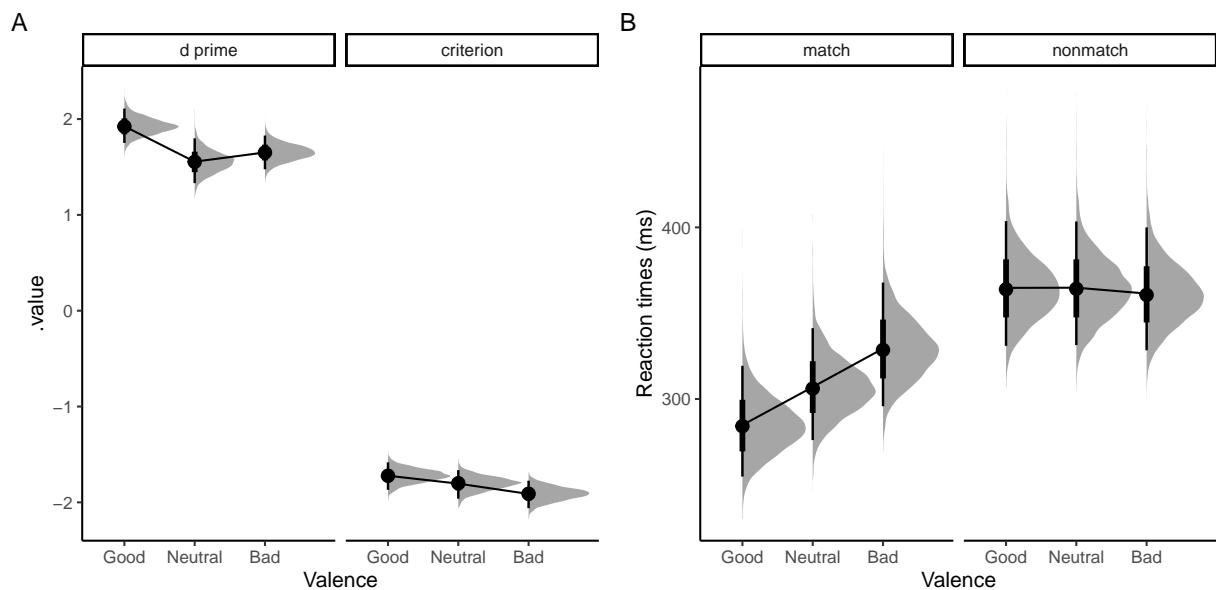


Figure 14. Exp6a: Results of Bayesian GLM analysis.

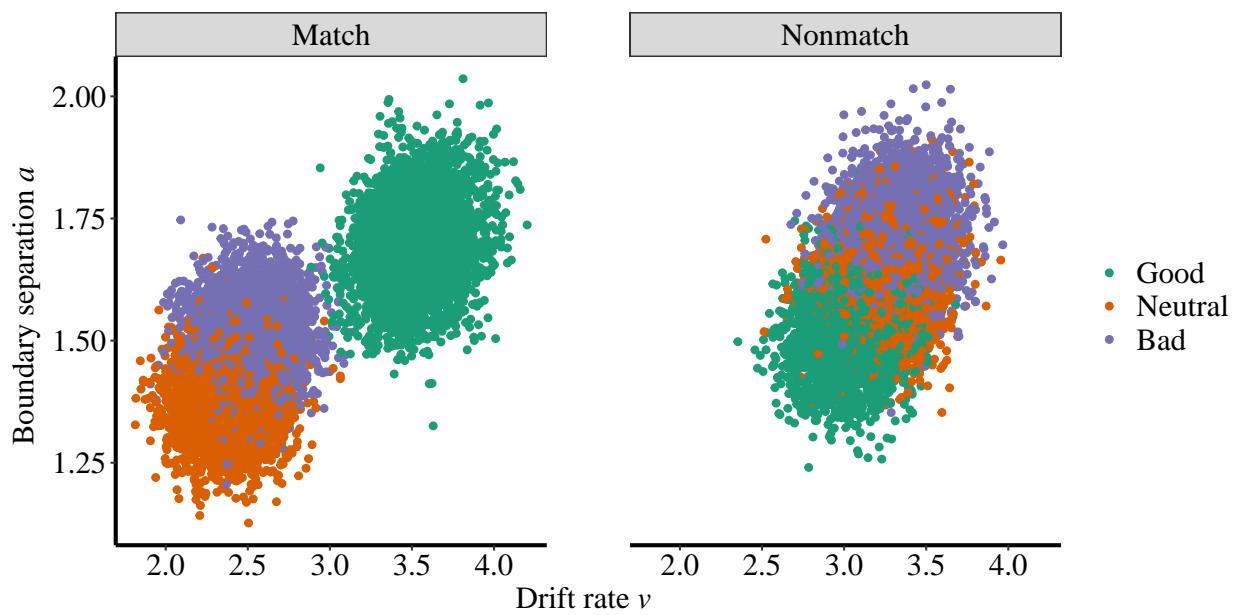


Figure 15. exp6a: Results of HDDM.

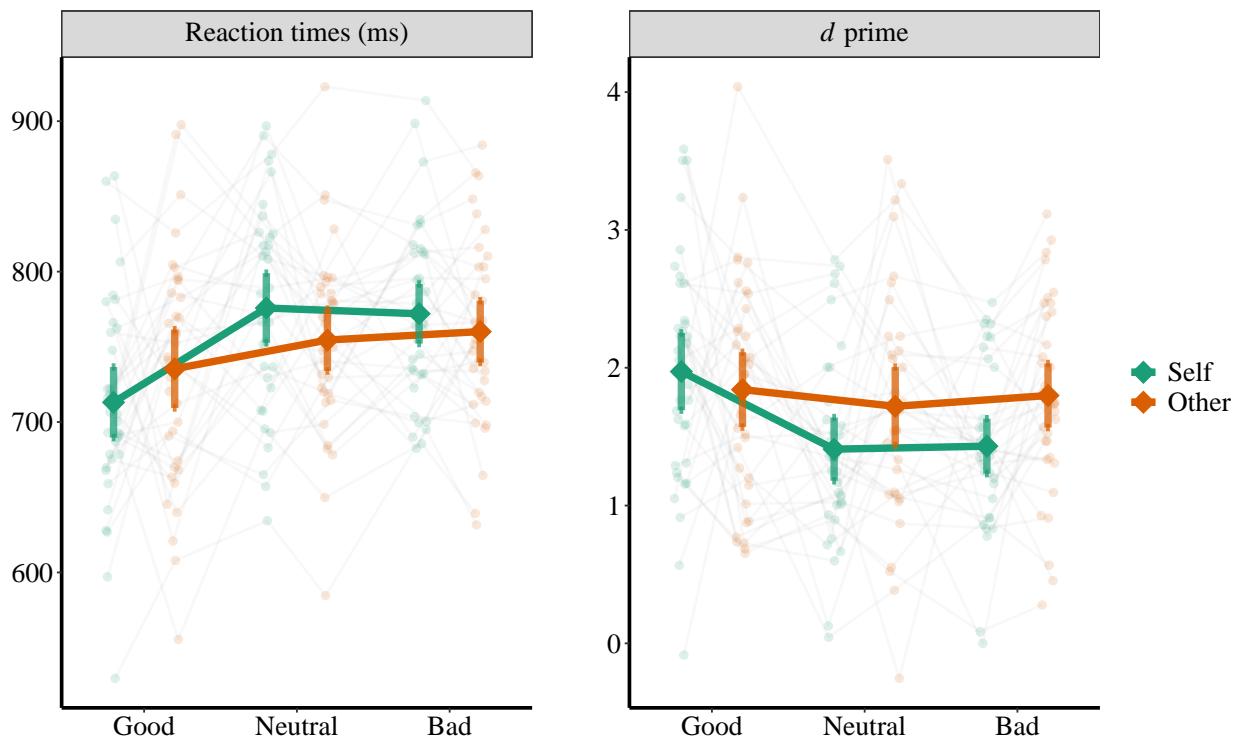


Figure 16. RT and d prime of Experiment 3a.

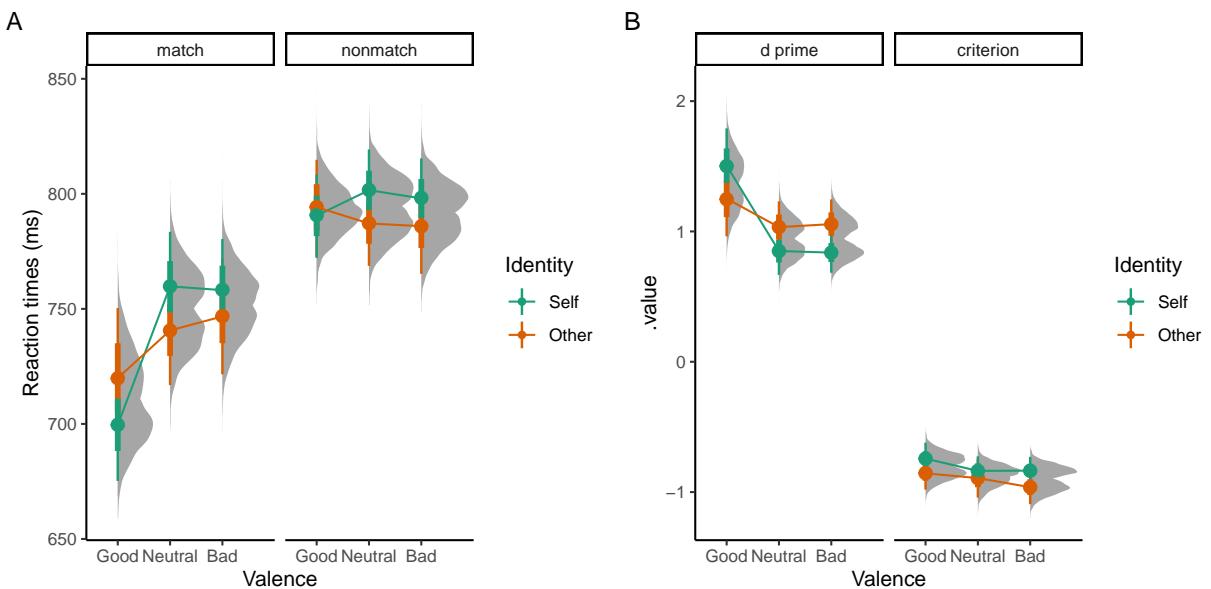


Figure 17. Exp3a: Results of Bayesian GLM analysis.

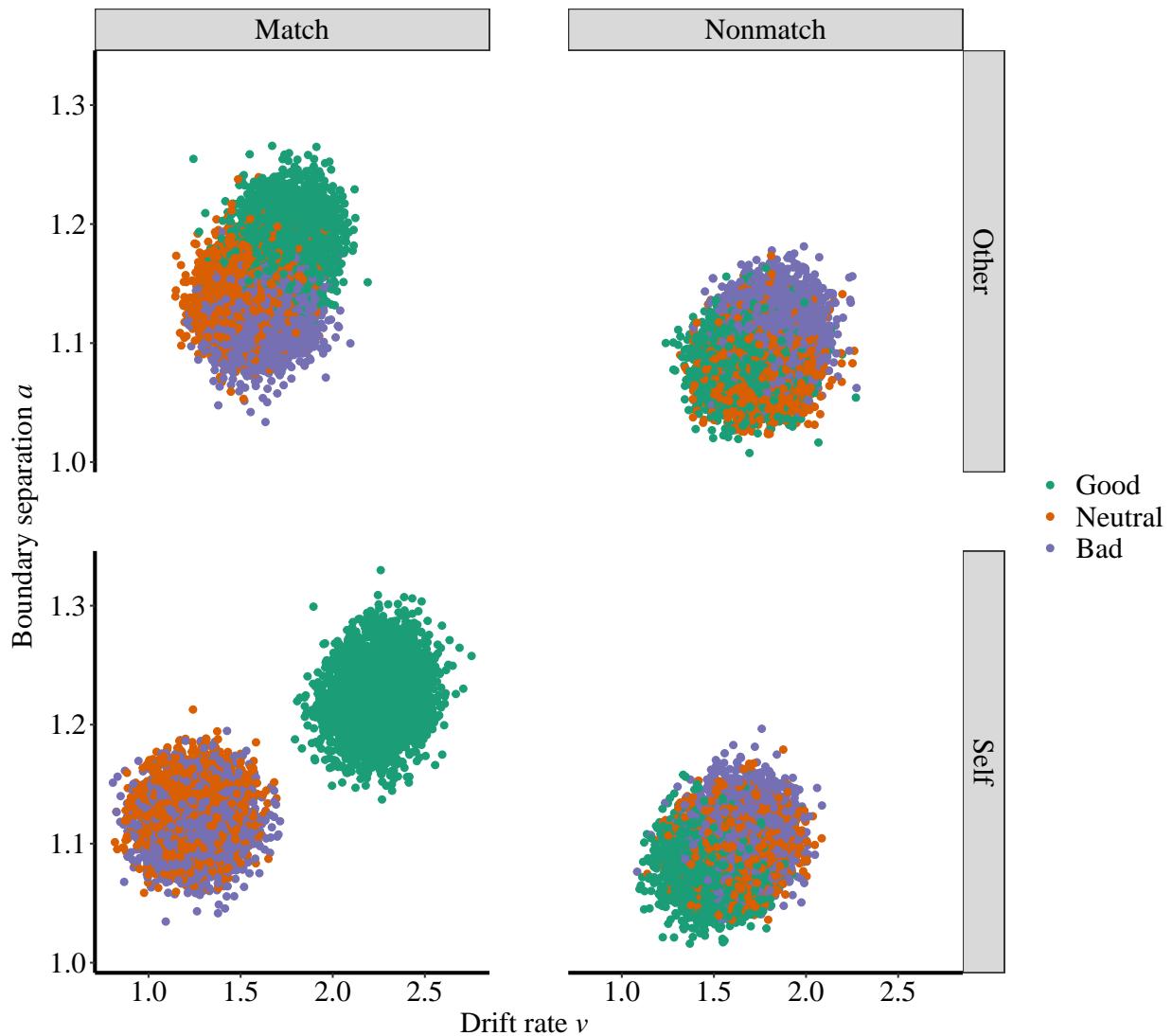


Figure 18. Exp3a: Results of HDDM.

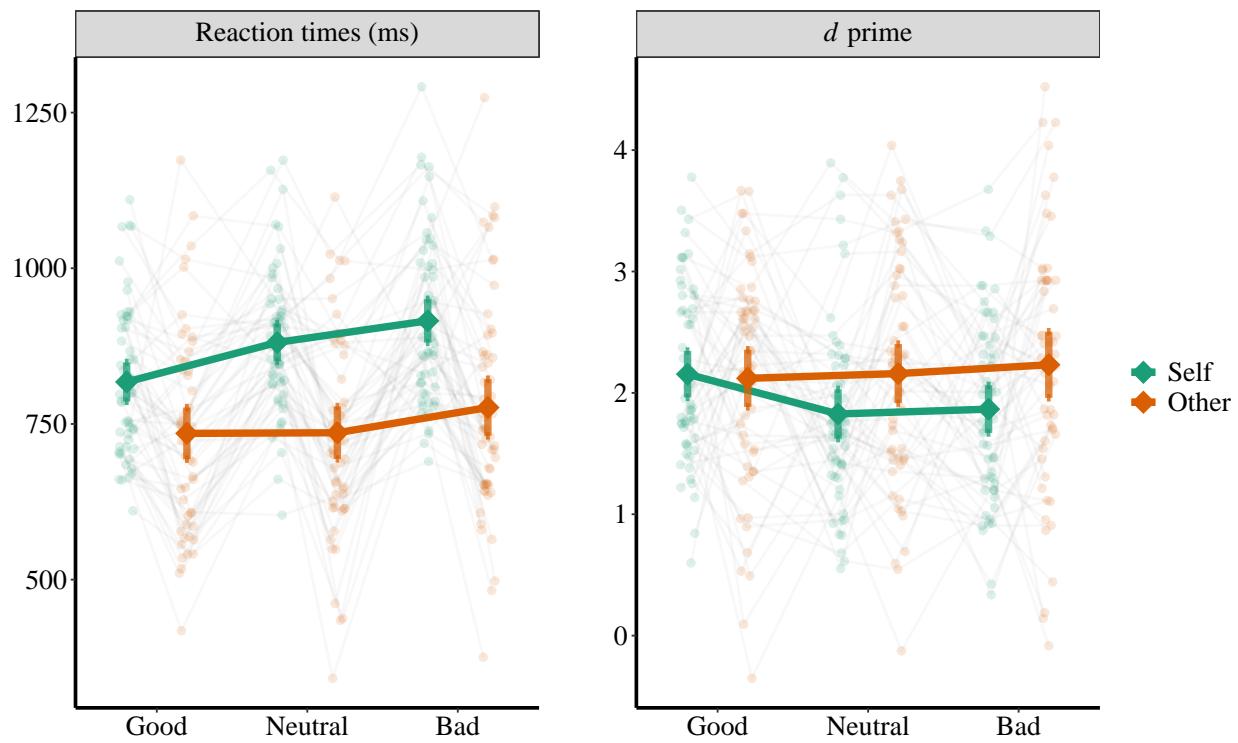


Figure 19. RT and d prime of Experiment 3b.

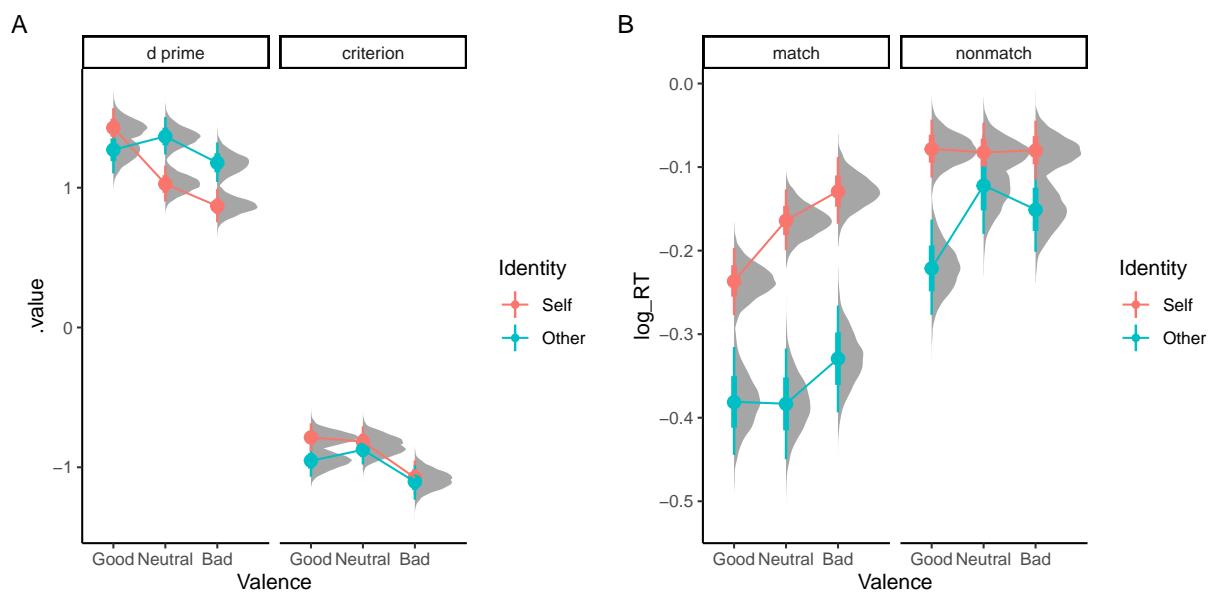


Figure 20. exp3b: Results of Bayesian GLM analysis.

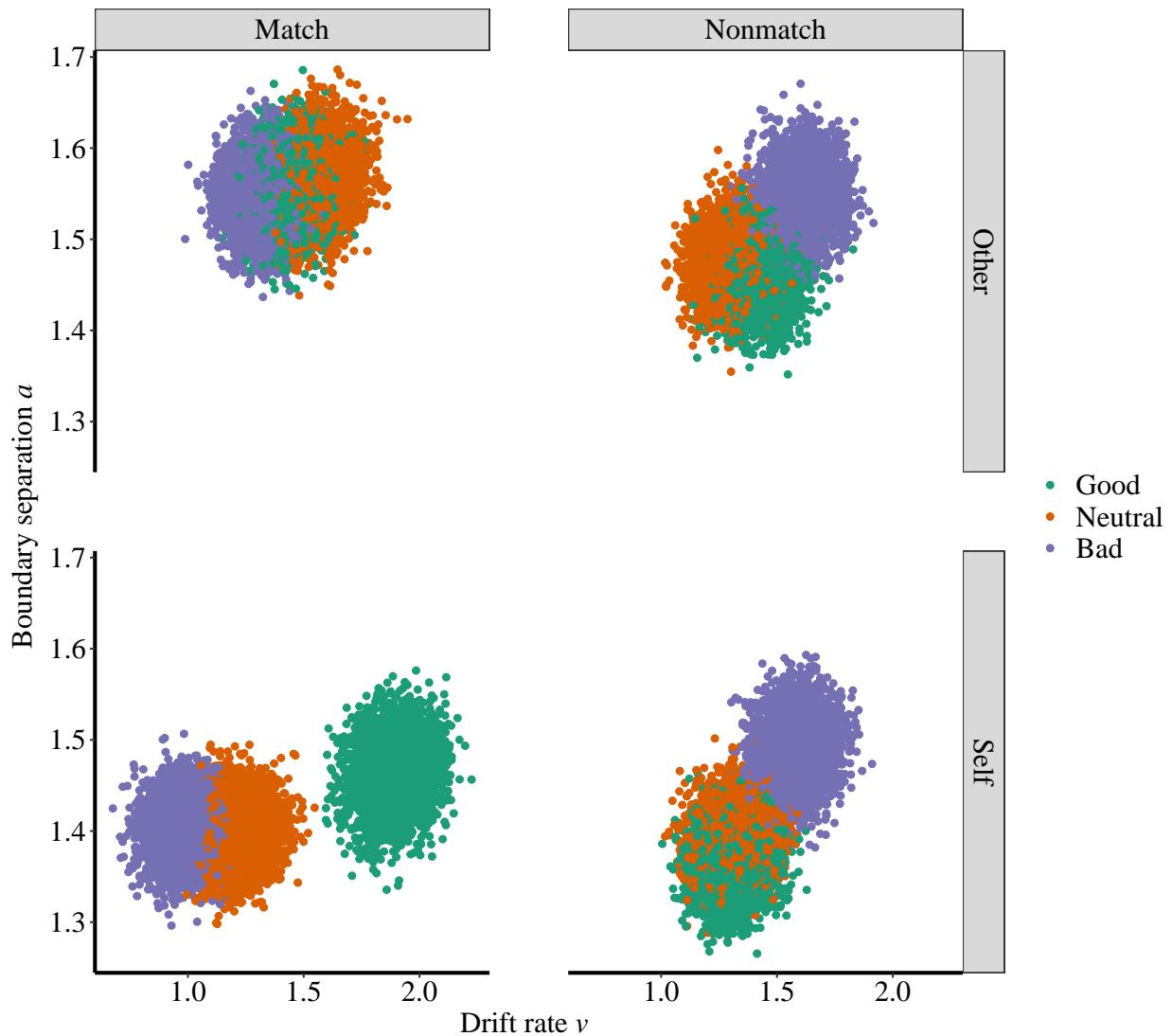


Figure 21. exp3b: Results of HDDM.

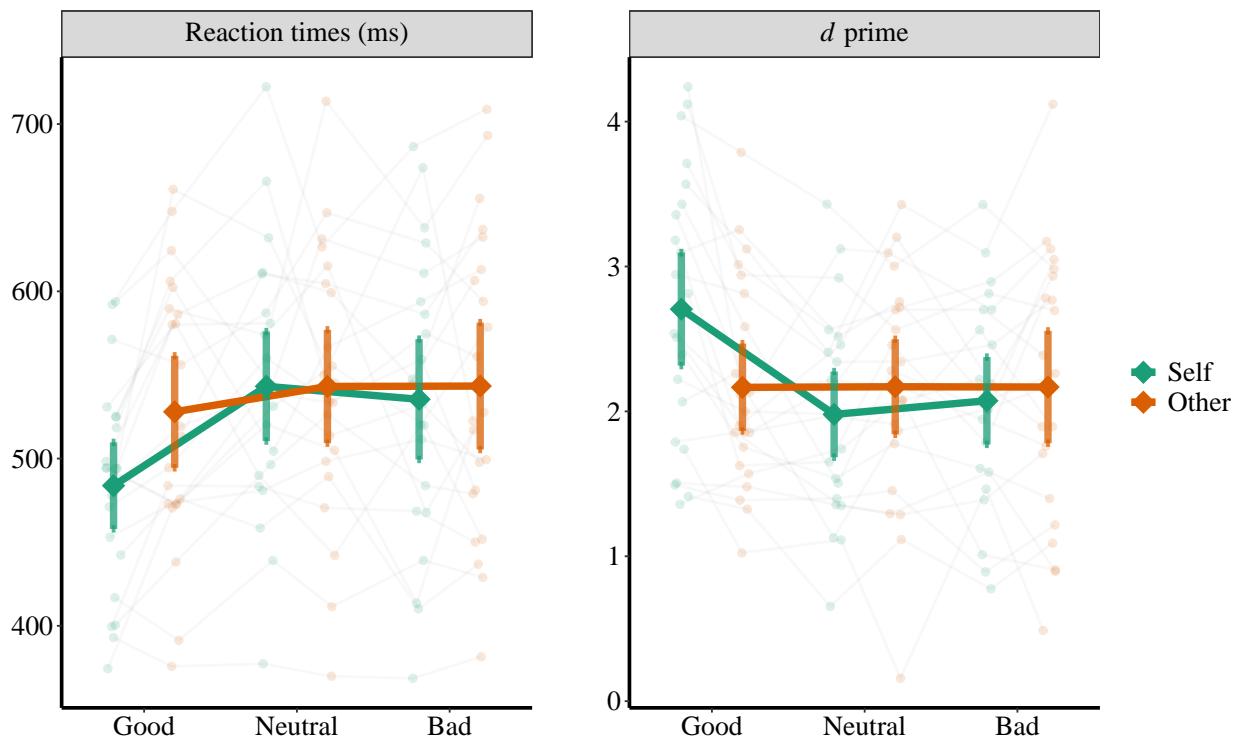


Figure 22. RT and d' prime of Experiment 6b.

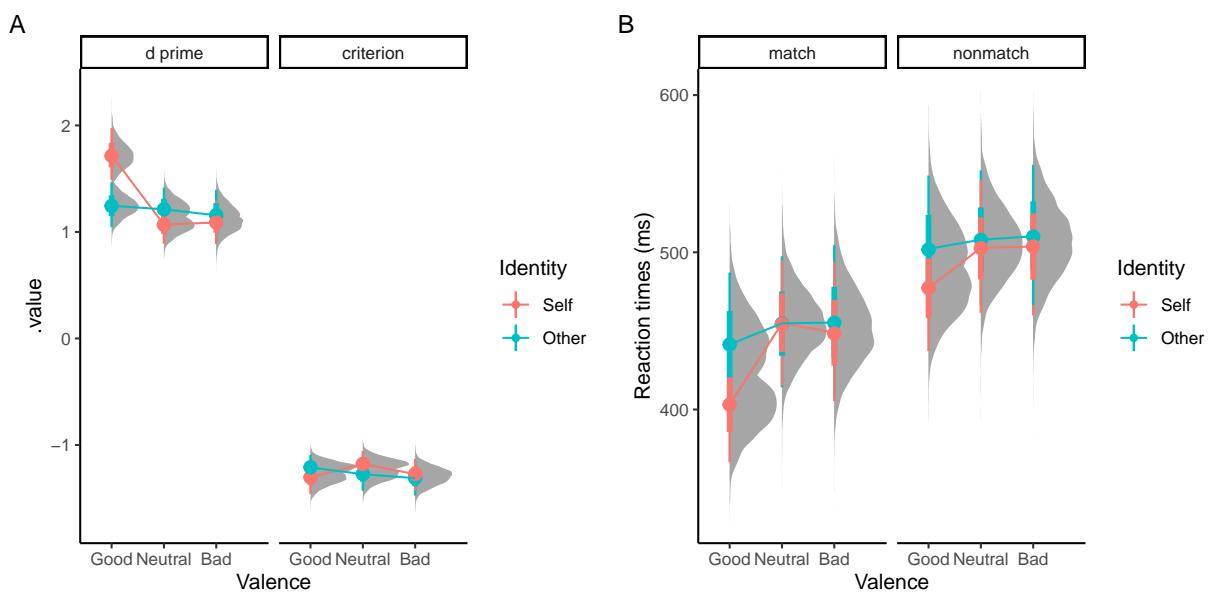


Figure 23. exp6b_d1: Results of Bayesian GLM analysis.

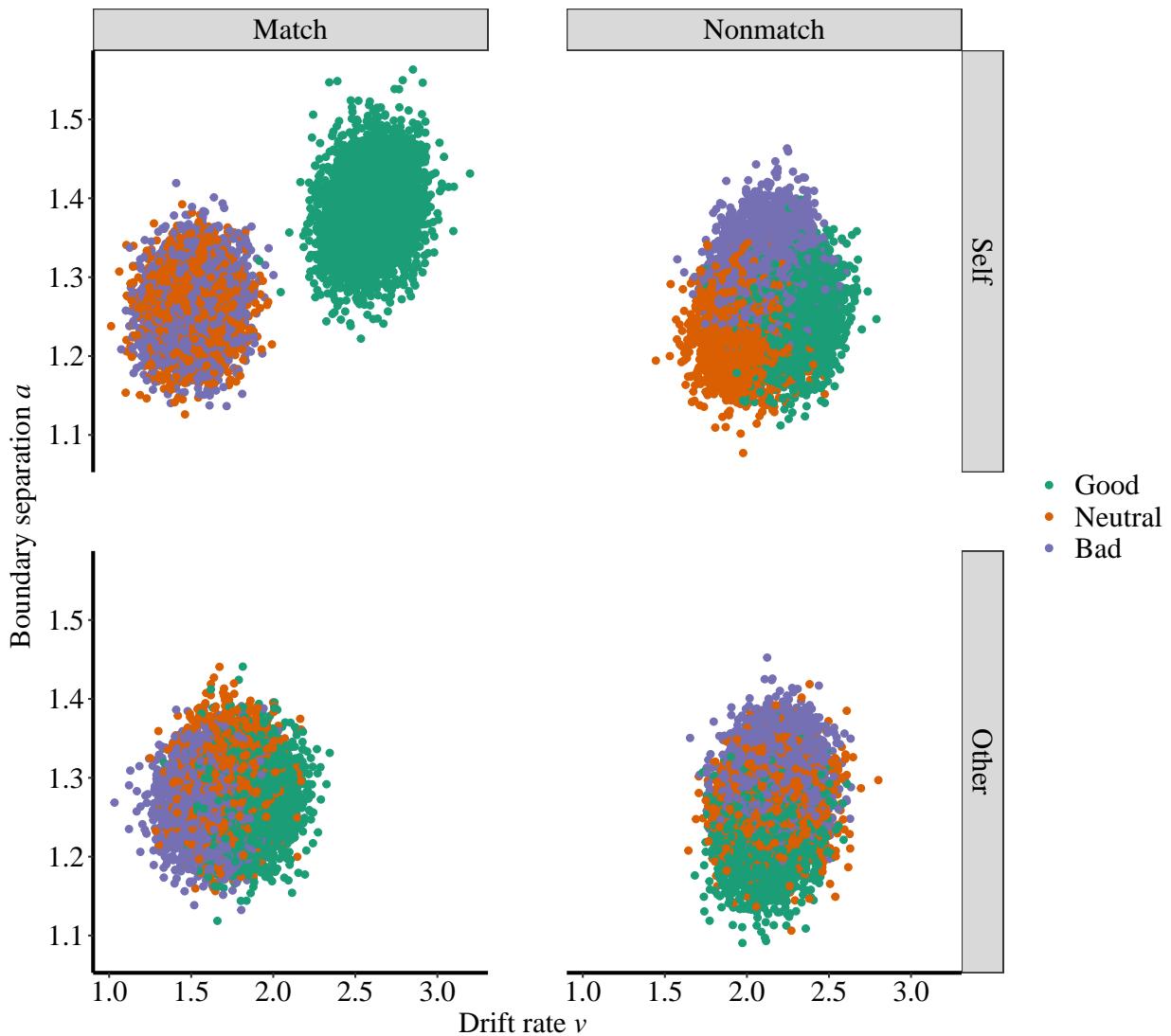


Figure 24. exp6b: Results of HDDM (Day 1).

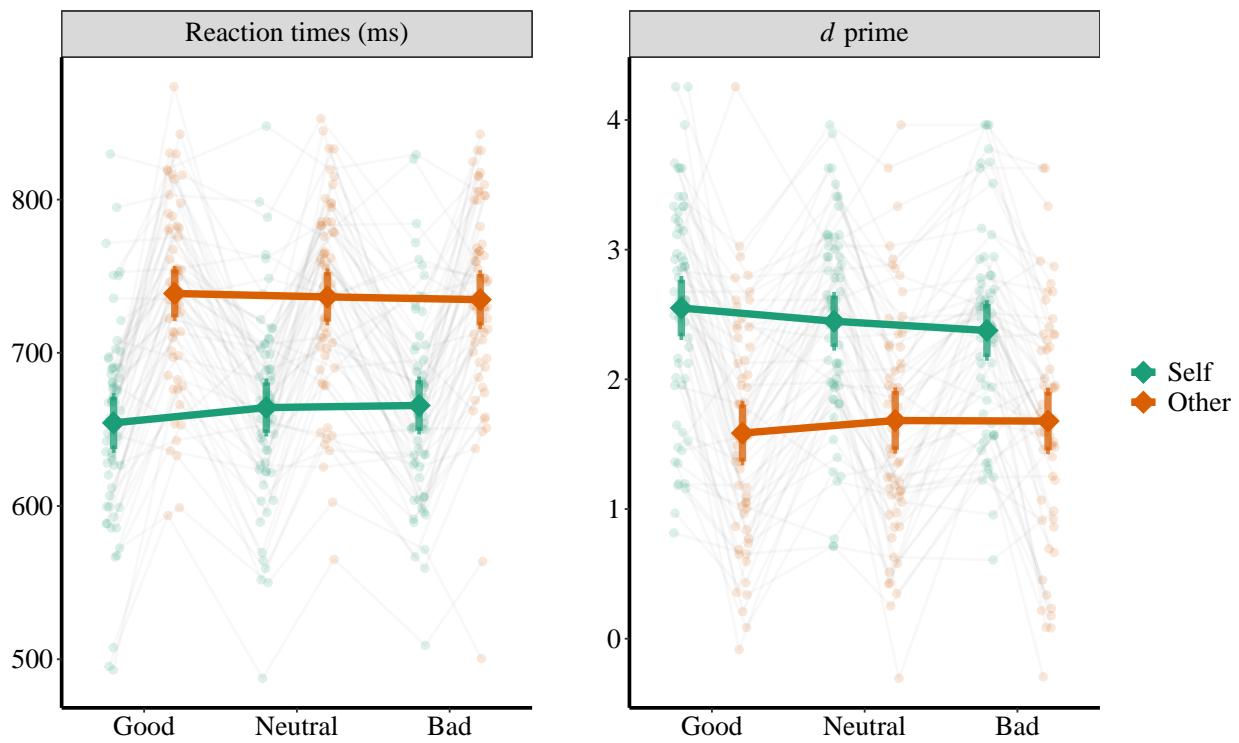


Figure 25. RT and d' of Experiment 4a.

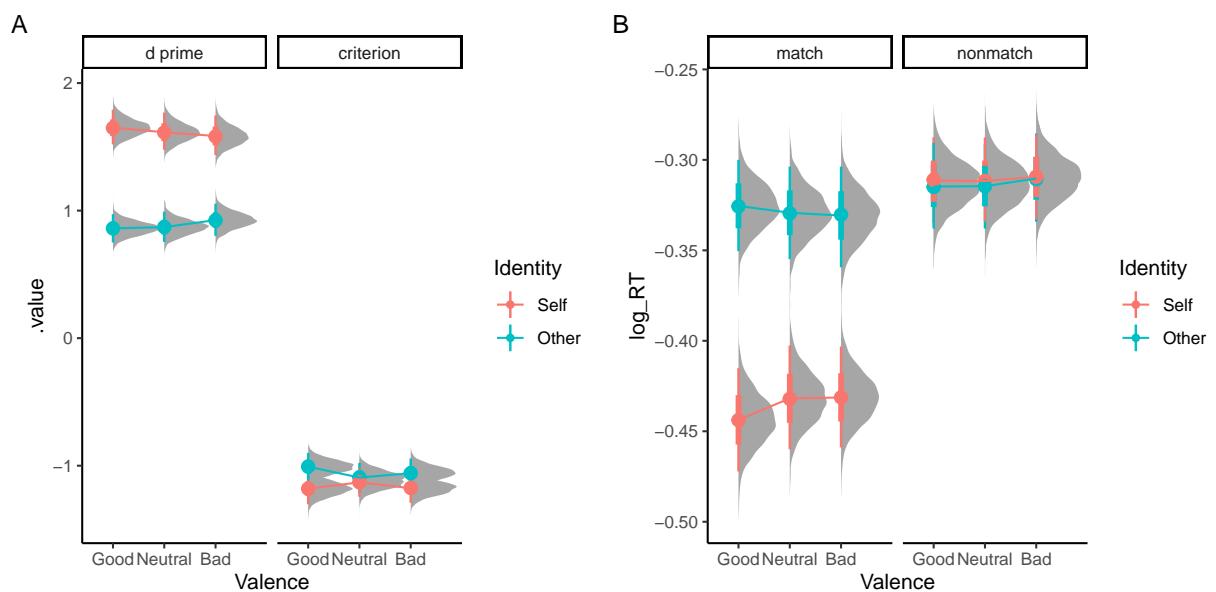


Figure 26. exp4a: Results of Bayesian GLM analysis.

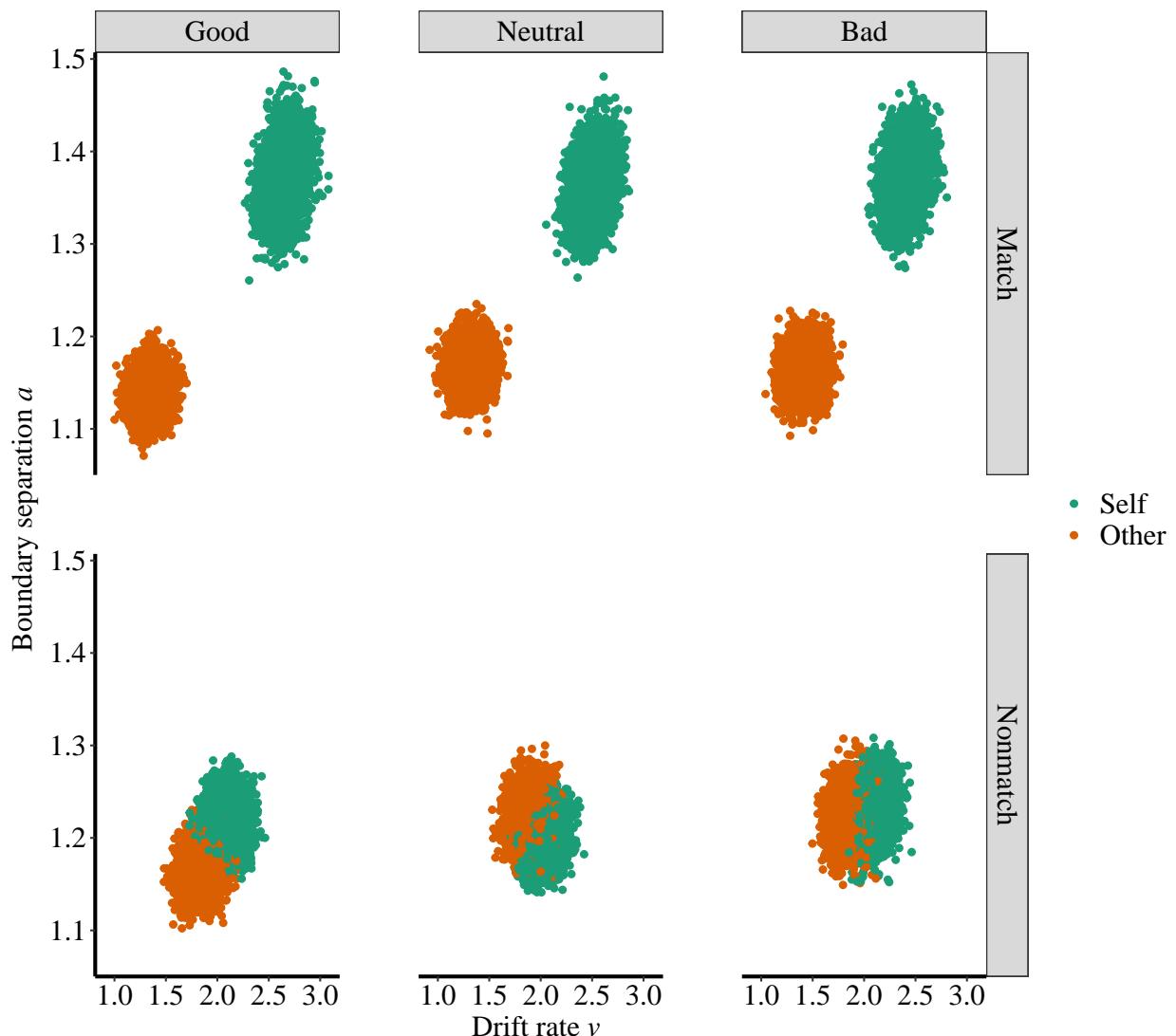


Figure 27. exp4a: Results of HDDM.

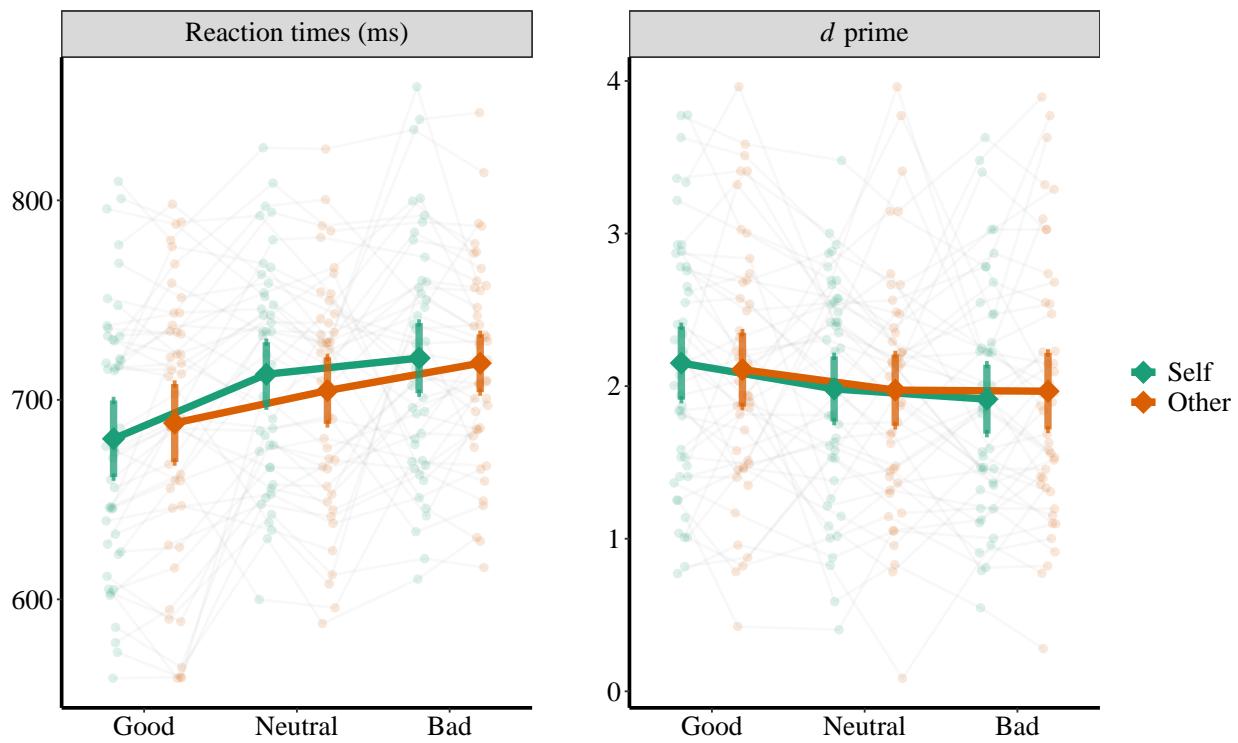


Figure 28. RT and d' prime of Experiment 4b.

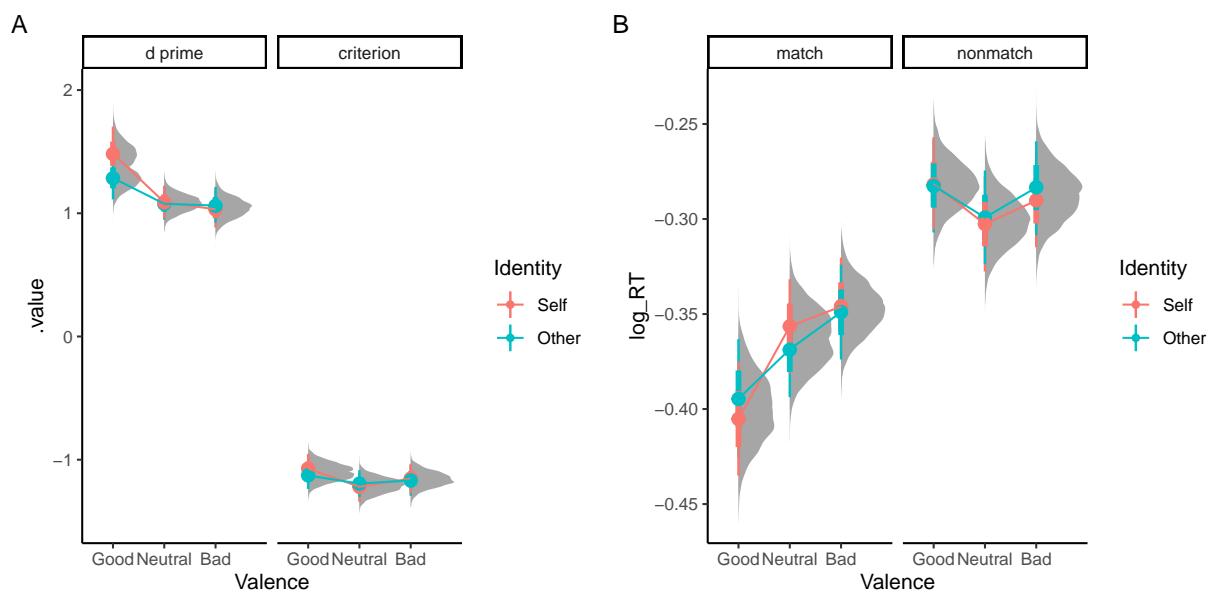


Figure 29. exp4b: Results of Bayesian GLM analysis.

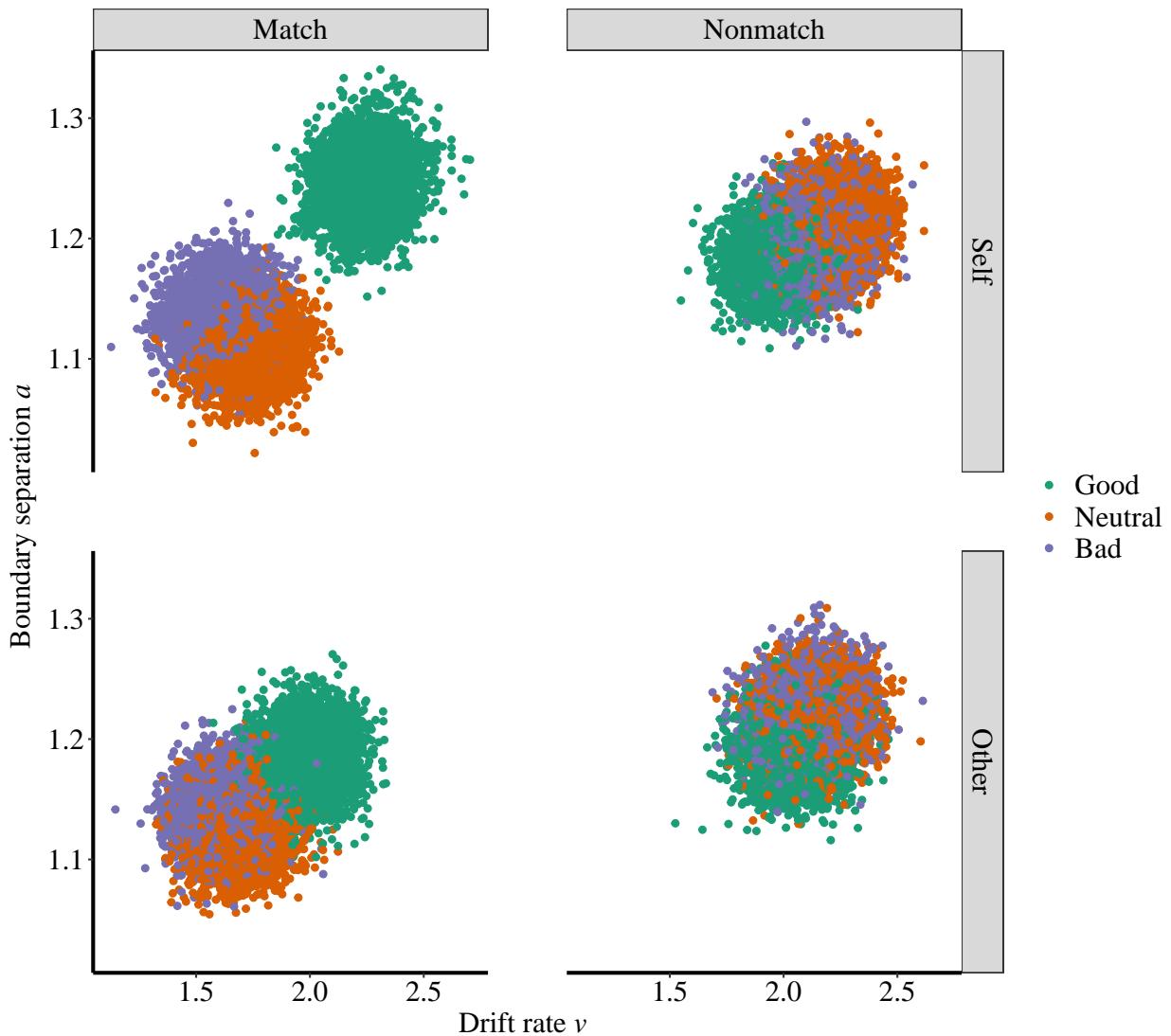


Figure 30. exp4b: Results of HDDM.

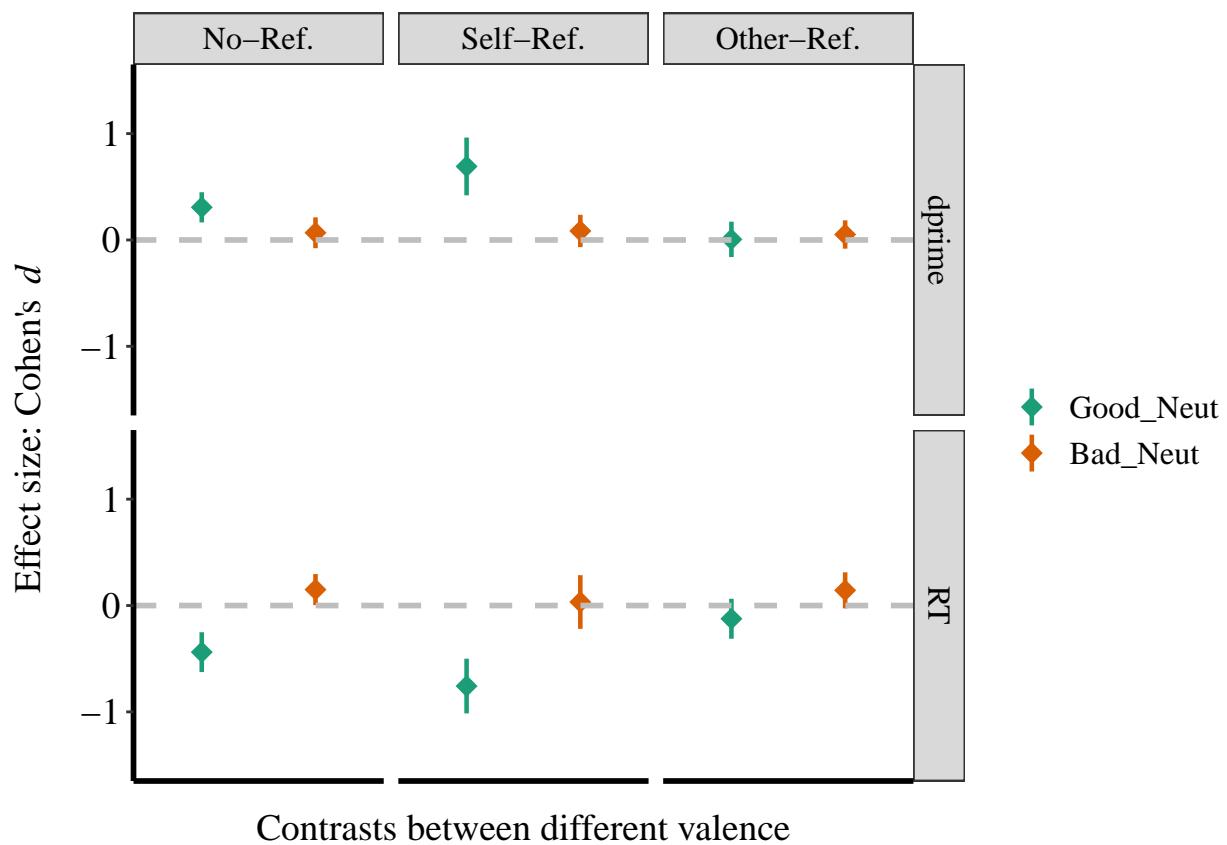


Figure 31. Effect size (Cohen's d) of Valence.

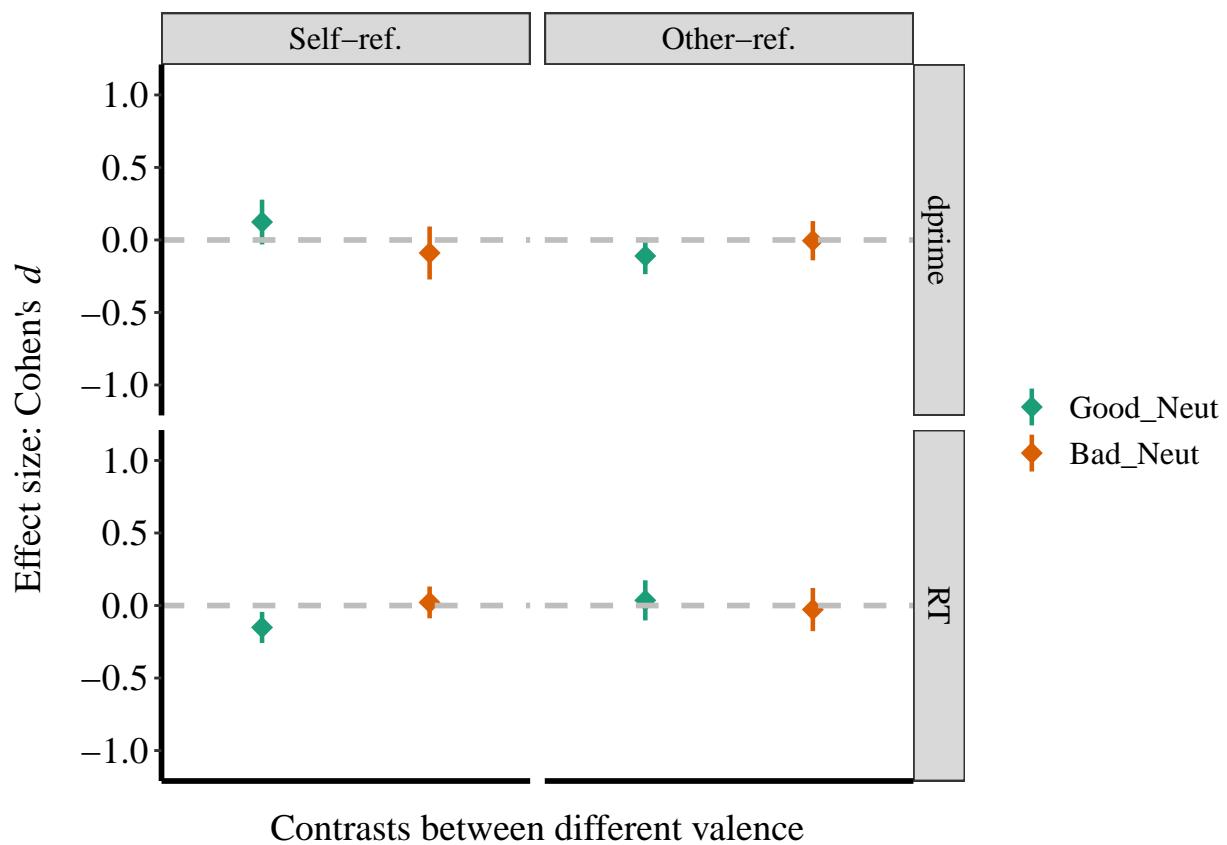


Figure 32. Effect size (Cohen's d) of Valence in Exp4a.

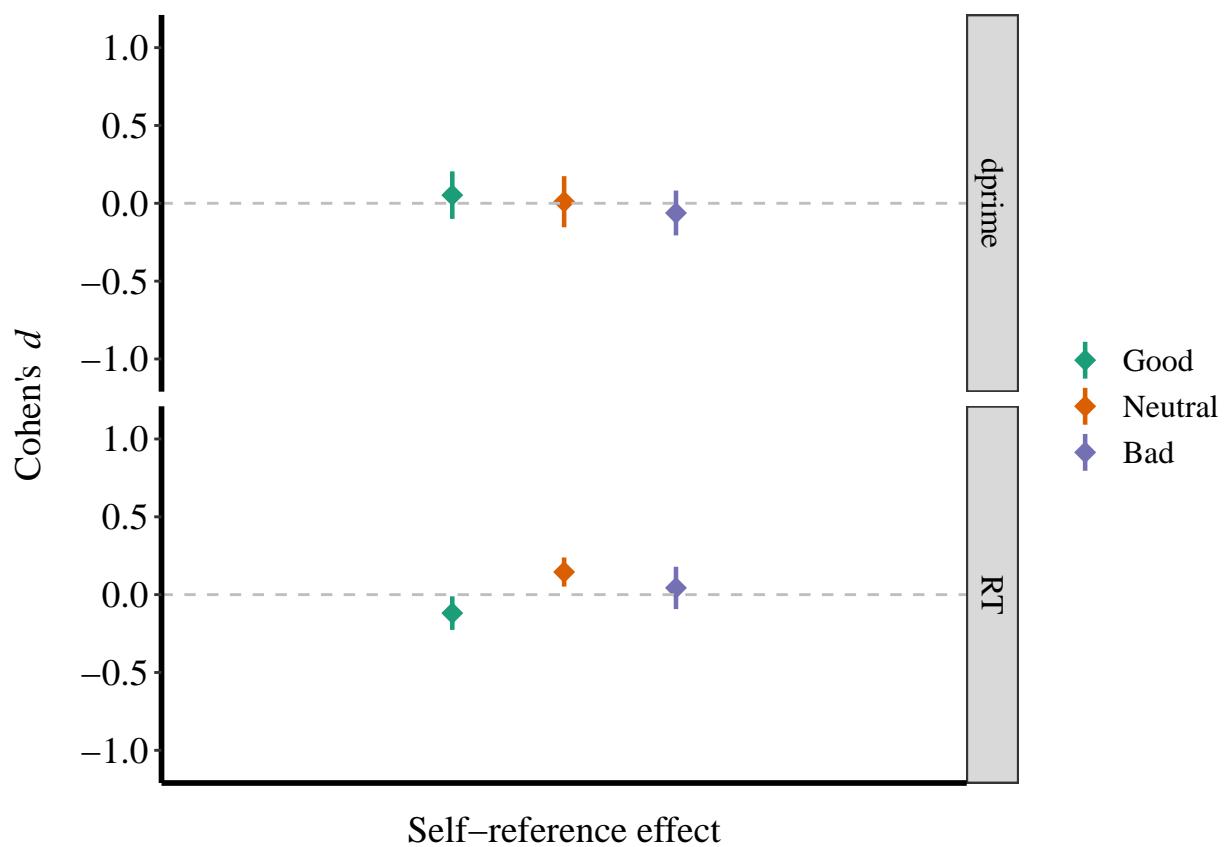


Figure 33. Effect size (Cohen's d) of Valence in Exp4b.

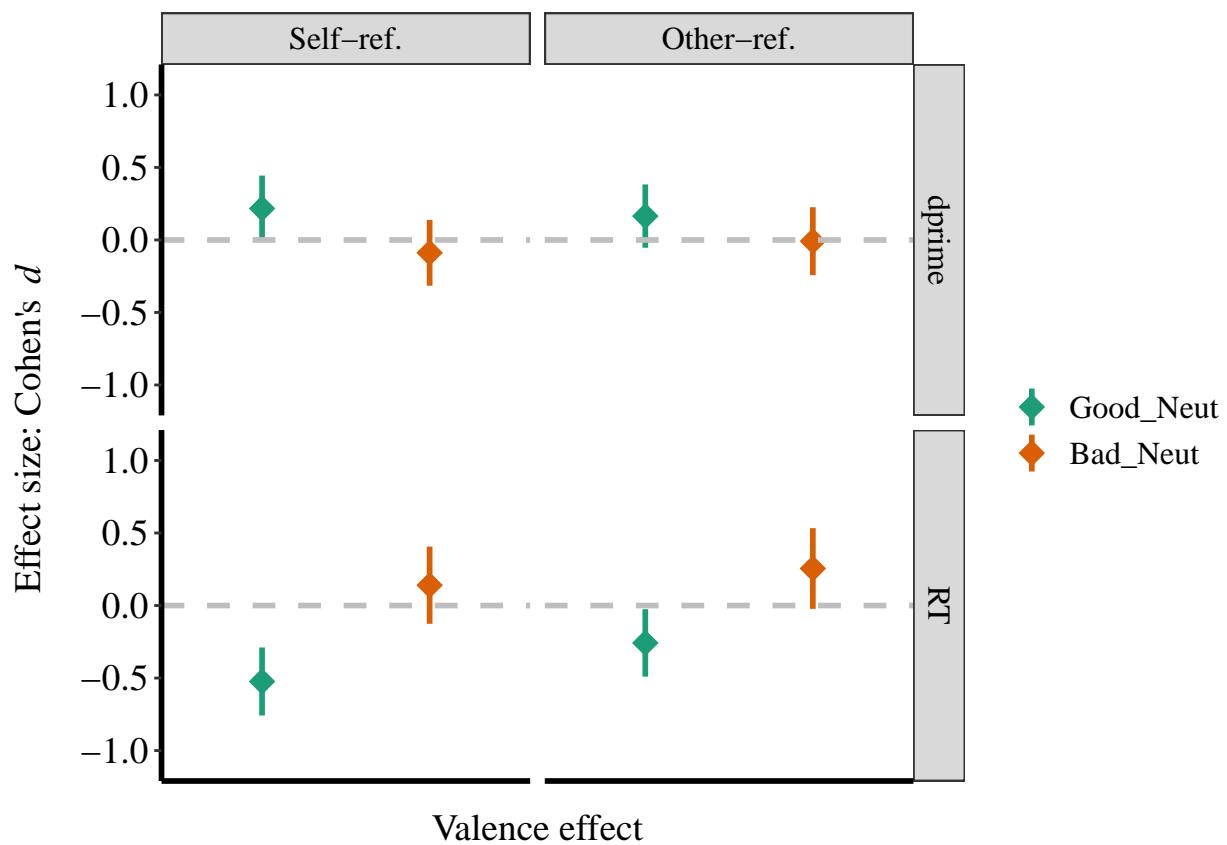


Figure 34. Effect size (Cohen's d) of Valence in Exp4b.

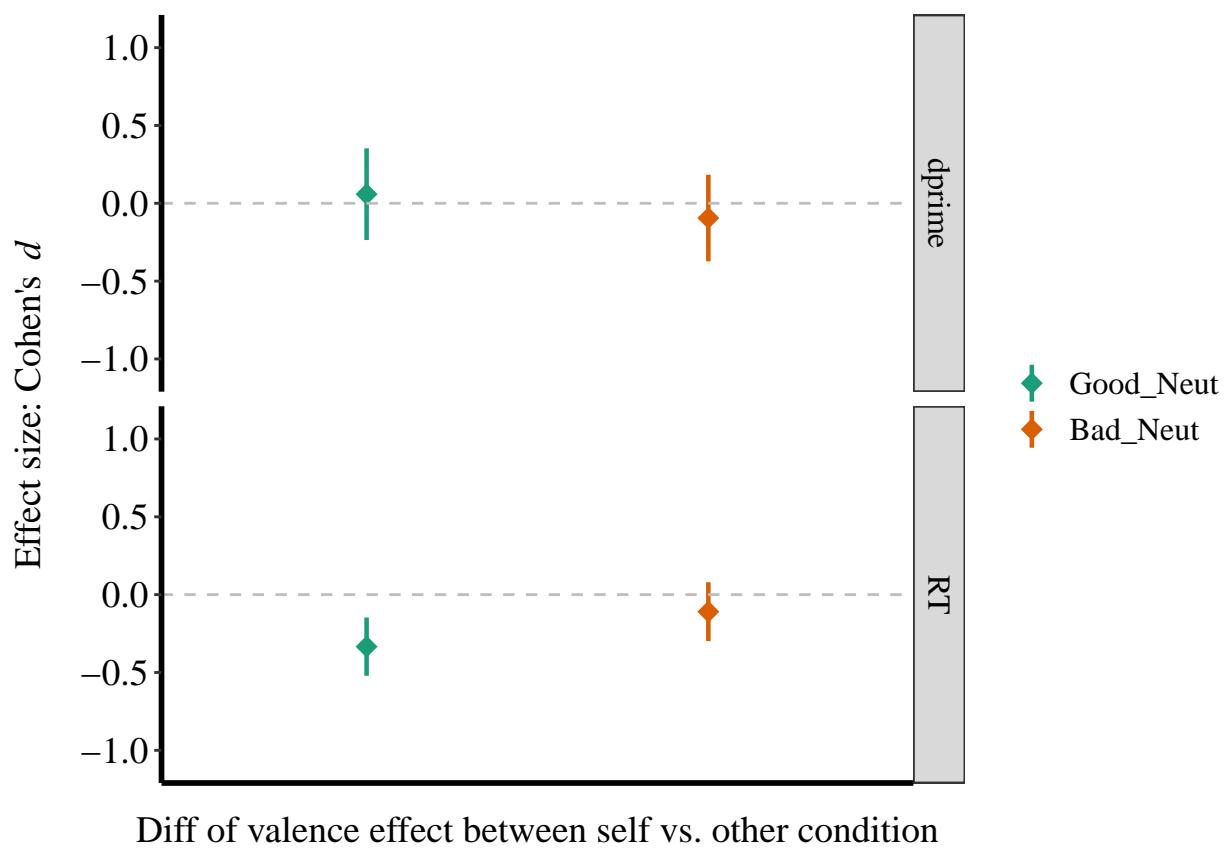


Figure 35. Effect size (Cohen's d) of Valence in Exp4b.

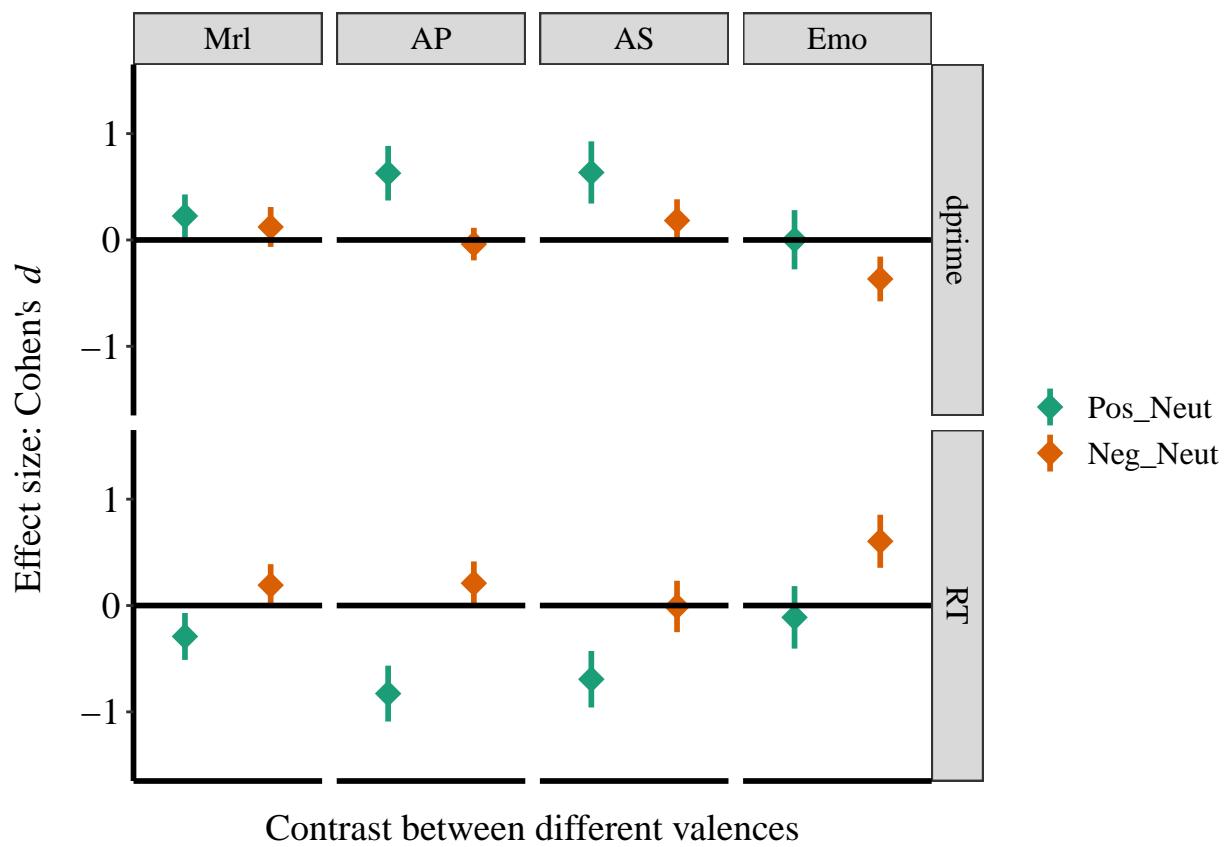


Figure 36. Effect size (Cohen's d) of Valence in Exp5.

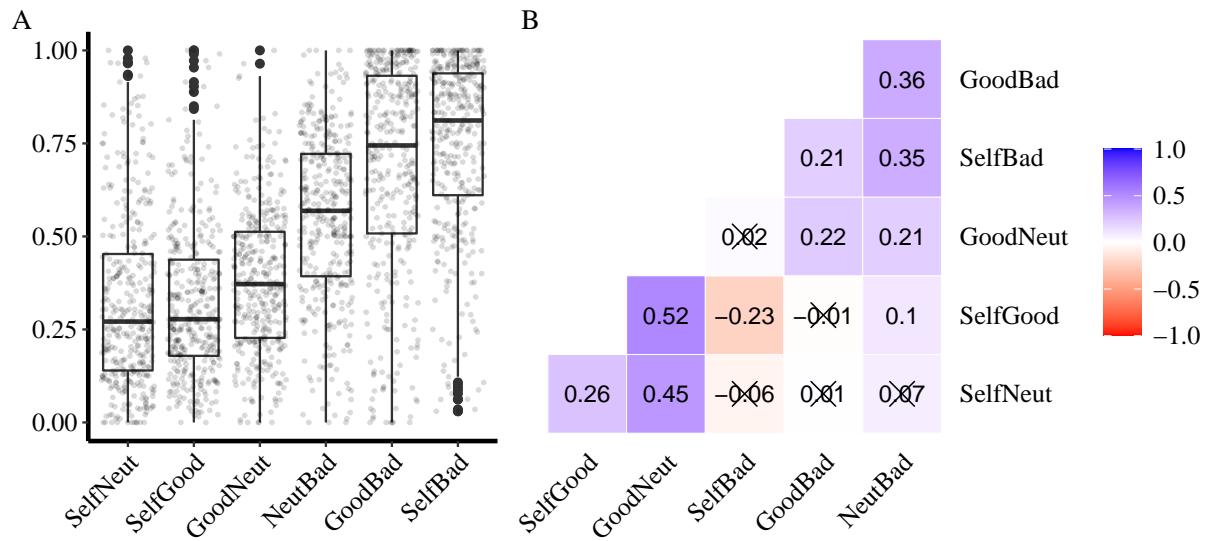


Figure 37. Self-rated personal distance

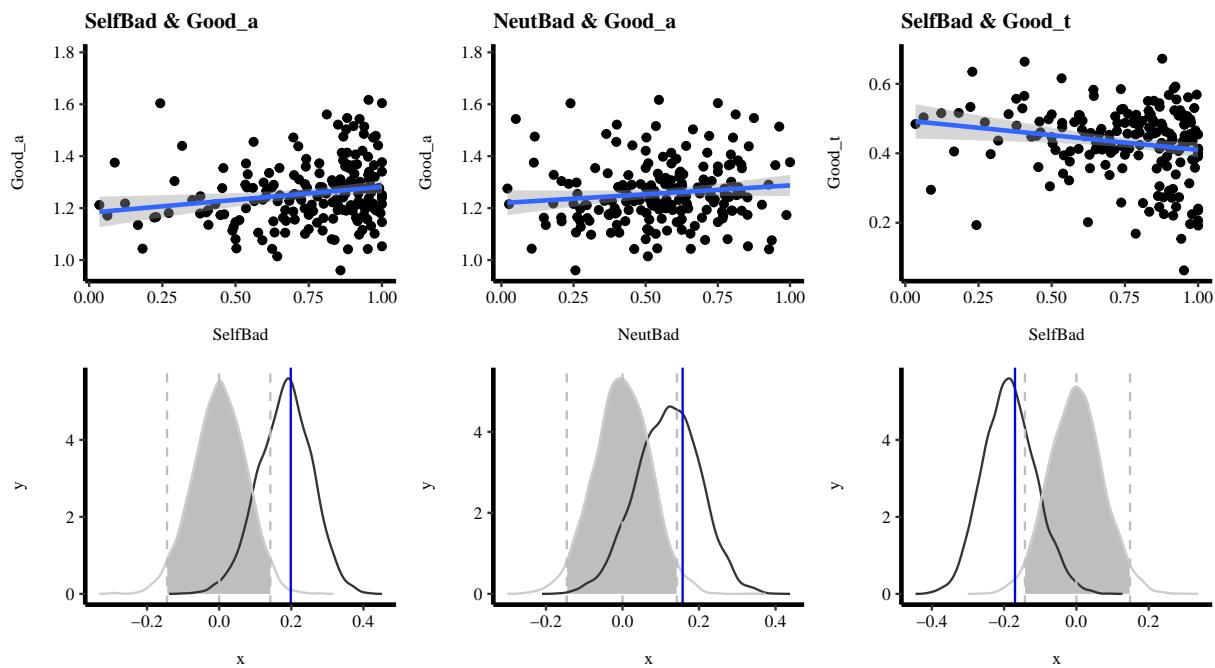


Figure 38. Correlation between moral identity and boundary separation of good condition; moral self-image and drift rate of good condition

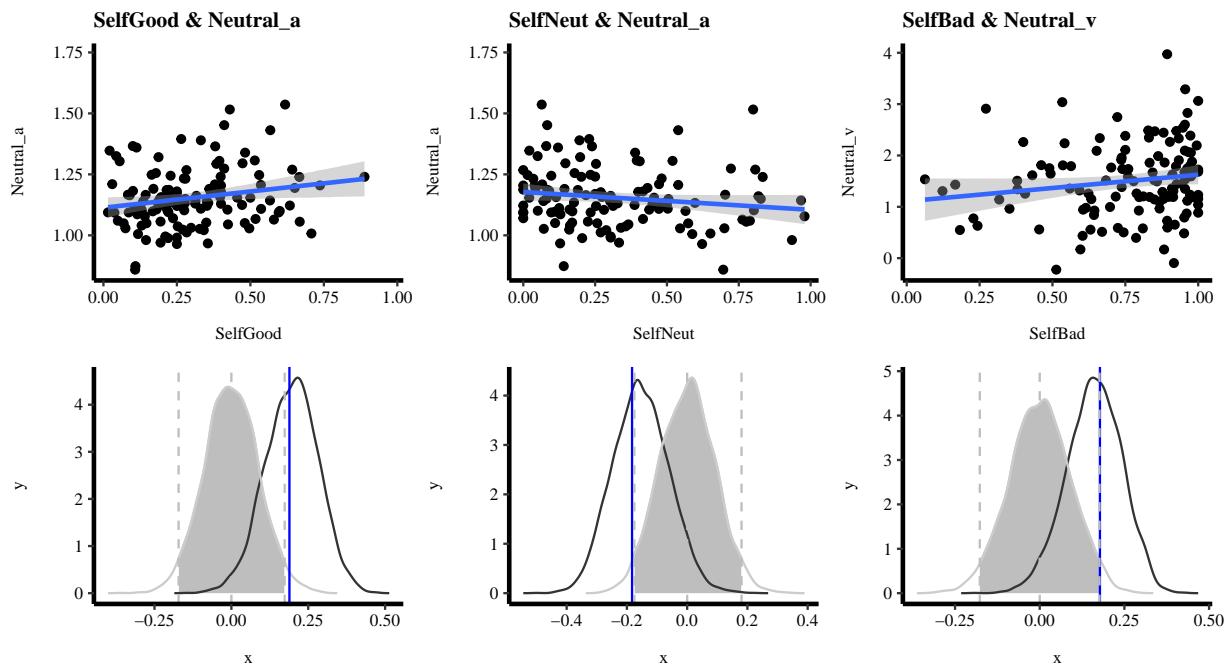


Figure 39. Correlation between personal distance and boundary separation of neutral condition