

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

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

³ ¹ TBA

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

⁵ ³ Tsinghua University, 100084 Beijing, China

⁶ ⁴ University of Aberdeen, Aberdeen, Scotland

⁷ Author Note

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

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

¹⁰ Psychology, University of Aberdeen, Aberdeen, Scotland.

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

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

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

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

¹⁵ Germany. E-mail: hcp4715@gmail.com

16

Abstract

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

37

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

38

Word count: X

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

40 # Introduction

41 Perceptual decision-making is an important window for understanding the cognition
42 (Shadlen & Kiani, 2013).

43 The role of perception in social cognition has been largely ignored, but accumulating
44 evidence revealed that perception provides rich information about our social cognition
45 (Stolier & Freeman, 2016; Xiao, Coppin, & Bavel, 2016).

46 Here we investigated how the instantly learned moral valence changed the perceptual
47 decision-making and the underlying psychological processes.

48 Given the importance of morality in social life (DeScioli, 2016) and identity (Freitas,
49 Cikara, Grossmann, & Schlegel, 2017; Strohminger, Knobe, & Newman, 2017; Zhang,
50 Chen, Schlegel, Hicks, & Chen, 2019), and evidence that moral character impacts how
51 people evaluate themselves [XXXX], desired personality change (Sun & Goodwin, 2020),
52 memory (Carlson, Maréchal, Oud, Fehr, & Crockett, 2020; Kouchaki & Gino, 2016; Shu,
53 Gino, & Bazerman, 2011; Stanley & De Brigard, 2019), one would expect that moral
54 character, the morality related trait, will also reflected in perceptual decision-making. Yet,
55 this effect was less studied (Anderson et al., 2011; Gantman & Van Bavel, 2014). This
56 moral perception effect was driven by motivation (Gantman & Van Bavel, 2016) and
57 influenced the information spreading online (Brady, Gantman, & Van Bavel, 2020).
58 Especially lacking is the process of this effect.

59 The current study first explored and confirmed a positive effect of moral character in
60 perceptual decision-making, using an associative learning task, then attempted to provide a
61 mechanistic explanation for this positivity effect: spontaneous self-identification with the
62 moral good character (Juechems and Summerfield (2019): Self-relevance is an important
63 dimension in determine the value, i.e., intrinsic goal), both implicitly and explicitly.

64 Finally, this positivity effect was also found in other social traits (beauty) but not
65 non-social, emotional states.

66 Potential theoretical discussion points: Close distance of the semantic representation
67 of self and moral character (attractor network) (Freeman & Ambady, 2011). The
68 core/true/authentic self concept.

69 We reported behavioral results from eleven experiments. In first set of experiments,
70 we found that shapes associated with morally positive person label were responded faster
71 and more accurately. In the second set of experiments, we explore the potential role of
72 good self in perceptual matching task and added one more independent variable, we found
73 that the effect was mainly on good self. In the third part we tested whether the morality
74 will automatically binds with self but not other. Finally, we explore the correlation
75 between behavioral task and questionnaire scores.

76 **Disclosures**

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

81 All the experiments reported were not pre-registered. Most experiments (1a ~ 6b,
82 except experiment 3b) reported in the current study were first finished between 2014 to
83 2016 in Tsinghua University, Beijing, China. Participants in these experiments were
84 recruited in the local community. To increase the sample size of experiments to 50 or more
85 (Simmons, Nelson, & Simonsohn, 2013), we recruited additional participants in Wenzhou
86 University, Wenzhou, China in 2017 for experiment 1a, 1b, 4a, and 4b. Experiment 3b was
87 finished in Wenzhou University in 2017. To have a better estimation of the effect size, we
88 included the data from two experiments (experiment 7a, 7b) that were reported in Hu,

⁸⁹ Lan, Macrae, and Sui (2020) (See Table S1 for overview of these experiments).

⁹⁰ All participant received informed consent and compensated for their time. These
⁹¹ experiments were approved by the ethic board in the Department of Tsinghua University.

⁹² **General methods**

⁹³ **Design and Procedure**

⁹⁴ This series of experiments started to test the effect of instantly acquired true self
⁹⁵ (moral self) on perceptual decision-making. For this purpose, we used the social associative
⁹⁶ learning paradigm (or tagging paradigm)(Sui, He, & Humphreys, 2012), in which
⁹⁷ participants first learned the associations between geometric shapes and labels of person
⁹⁸ with different moral character (e.g., in first three studies, the triangle, square, and circle
⁹⁹ and good person, neutral person, and bad person, respectively). The associations of the
¹⁰⁰ shapes and label were counterbalanced across participants. After remembered the
¹⁰¹ associations, participants finished a practice phase to familiar with the task, in which they
¹⁰² viewed one of the shapes upon the fixation while one of the labels below the fixation and
¹⁰³ judged whether the shape and the label matched the association they learned. When
¹⁰⁴ participants reached 60% or higher accuracy at the end of the practicing session, they
¹⁰⁵ started the experimental task which was the same as in the practice phase.

¹⁰⁶ The experiment 1a, 1b, 1c, 2, and 6a shared a 2 (matching: match vs. nonmatch) by
¹⁰⁷ 3 (moral valence: good vs. neutral vs. bad) within-subject design. Experiment 1a was the
¹⁰⁸ first one of the whole series studies and 1b, 1c, and 2 were conducted to exclude the
¹⁰⁹ potential confounding factors. More specifically, experiment 1b used different Chinese
¹¹⁰ words as label to test whether the effect only occurred with certain familiar words.
¹¹¹ Experiment 1c manipulated the moral valence indirectly: participants first learned to
¹¹² associate different moral behaviors with different neutral names, after remembered the
¹¹³ association, they then performed the perceptual matching task by associating names with

114 different shapes. Experiment 2 further tested whether the way we presented the stimuli
115 influence the effect of valence, by sequentially presenting labels and shapes. Note that part
116 of participants of experiment 2 were from experiment 1a because we originally planned a
117 cross task comparison. Experiment 6a, which shared the same design as experiment 2, was
118 an EEG experiment which aimed at exploring the neural correlates of the effect. But we
119 will focus on the behavioral results of experiment 6a in the current manuscript.

120 For experiment 3a, 3b, 4a, 4b, 6b, 7a, and 7b, we included self-reference as another
121 within-subject variable in the experimental design. For example, the experiment 3a directly
122 extend the design of experiment 1a into a 2 (matchness: match vs. nonmatch) by 2
123 (reference: self vs. other) by 3 (moral valence: good vs. neutral vs. bad) within-subject
124 design. Thus in experiment 3a, there were six conditions (good-self, neutral-self, bad-self,
125 good-other, neutral-other, and bad-other) and six shapes (triangle, square, circle, diamond,
126 pentagon, and trapezoids). The experiment 6b was an EEG experiment extended from
127 experiment 3a but presented the label and shape sequentially. Because of the relatively
128 high working memory load (six label-shape pairs), experiment 6b were conducted in two
129 days: the first day participants finished perceptual matching task as a practice, and the
130 second day, they finished the task again while the EEG signals were recorded. Experiment
131 3b was designed to separate the self-referential trials and other-referential trials. That is,
132 participants finished two different blocks: in the self-referential blocks, they only responded
133 to good-self, neutral-self, and bad-self, with half match trials and half non-match trials; for
134 the other-reference blocks, they only responded to good-other, neutral-other, and
135 bad-other. Experiment 7a and 7b were designed to test the cross task robustness of the
136 effect we observed in the aforementioned experiments (see, Hu et al., 2020). The matching
137 task in these two experiments shared the same design with experiment 3a, but only with
138 two moral valence, i.e., good vs. bad. We didn't include the neutral condition in
139 experiment 7a and 7b because we found that the neutral and bad conditions constantly
140 showed non-significant results in experiment 1 ~ 6.

141 Experiment 4a and 4b were design to test the automaticity of the binding between
142 self/other and moral valence. In 4a, we used only two labels (self vs. other) and two shapes
143 (circle, square). To manipulate the moral valence, we added the moral-related words within
144 the shape and instructed participants to ignore the words in the shape during the task. In
145 4b, we reversed the role of self-reference and valence in the task: participant learnt three
146 labels (good-person, neutral-person, and bad-person) and three shapes (circle, square, and
147 triangle), and the words related to identity, “self” or “other”, were presented in the shapes.
148 As in 4a, participants were told to ignore the words inside the shape during the task.

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

155 E-prime 2.0 was used for presenting stimuli and collecting behavioral responses,
156 except that experiment 7a and 7b used Matlab Psychtoolbox (Brainard, 1997; Pelli, 1997).
157 For participants recruited in Tsinghua University, they finished the experiment individually
158 in a dim-lighted chamber, stimuli were presented on 22-inch CRT monitors and their head
159 were fixed by a chin-rest brace. The distance between participants' eyes and the screen was
160 about 60 cm. The visual angle of geometric shapes was about $3.7^\circ \times 3.7^\circ$, the fixation cross
161 is of ($0.8^\circ \times 0.8^\circ$ of visual angle) at the center of the screen. The words were of $3.6^\circ \times 1.6^\circ$
162 visual angle. The distance between the center of the shape or the word and the fixation
163 cross was 3.5° of visual angle. For participants recruited in Wenzhou University, they
164 finished the experiment in a group consisted of 3 ~ 12 participants in a dim-lighted testing
165 room. Participants were required to finished the whole experiment independently. Also,
166 they were instructed to start the experiment at the same time, so that the distraction
167 between participants were minimized. The stimuli were presented on 19-inch CRT monitor.

¹⁶⁸ The visual angles are could not be exactly controlled because participants's chin were not
¹⁶⁹ fixed.

¹⁷⁰ In most of these experiments, participant were also asked to fill a battery of
¹⁷¹ questionnaire after they finish the behavioral tasks. All the questionnaire data are open
¹⁷² (see, dataset 4 in Liu et al., 2020). See Table S1 for a summary information about all the
¹⁷³ experiments.

¹⁷⁴ **Data analysis**

¹⁷⁵ **Analysis of individual study.** We used the `tidyverse` of r (see script
¹⁷⁶ `Load_save_data.r`) to exclude the practicing trials, invalid trials of each participants, and
¹⁷⁷ invalid participants, if there were any, in the raw data. Results of each experiment were
¹⁷⁸ then analyzed in three different approaches.

¹⁷⁹ ***Classic NHST.***

¹⁸⁰ First, as in Sui et al. (2012), we analyzed the accuracy and reaction times using
¹⁸¹ classic repeated measures ANOVA in the Null Hypothesis Significance Test (NHST)
¹⁸² framework. Repeated measures ANOVAs is essentially a two-step mixed model. In the first
¹⁸³ step, we estimate the parameter on individual level, and in the second step, we used
¹⁸⁴ repeated ANOVA to test the Null hypothesis. More specifically, for the accuracy, we used a
¹⁸⁵ signal detection approach, in which individual' sensitivity d' was estimated first. To
¹⁸⁶ estimate the sensitivity, we treated the match condition as the signal while the nonmatch
¹⁸⁷ conditions as noise. Trials without response were coded either as "miss" (match trials) or
¹⁸⁸ "false alarm" (nonmatch trials). Given that the match and nonmatch trials are presented
¹⁸⁹ in the same way and had same number of trials across all studies, we assume that
¹⁹⁰ participants' inner distribution of these two types of trials had equal variance but may had
¹⁹¹ different means. That is, we used the equal variance Gaussian SDT model (EVSDT) here
¹⁹² (Rouder & Lu, 2005). The d' was then estimated as the difference of the standardized hit

¹⁹³ and false alarm rats (Stanislaw & Todorov, 1999):

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

¹⁹⁴ where the *HR* means hit rate and the *FAR* mean false alarm rate. *zHR* and *zFAR* are
¹⁹⁵ the standardized hit rate and false alarm rates, respectively. These two *z*-scores were
¹⁹⁶ converted from proportion (i.e., hit rate or false alarm rate) by inverse cumulative normal
¹⁹⁷ density function, Φ^{-1} (Φ is the cumulative normal density function, and is used convert *z*
¹⁹⁸ score into probabilities). Another parameter of signal detection theory, response criterion *c*,
¹⁹⁹ is defined by the negative standardized false alarm rate (DeCarlo, 1998): $-zFAR$.

²⁰⁰ For the reaction times (RTs), only RTs of accurate trials were analyzed. We first
²⁰¹ calculate the mean RTs of each participant and then subject the mean RTs of each
²⁰² participant to repeated measures ANOVA. Note that we set the alpha as .05. The repeated
²⁰³ measure ANOVA was done by *afex* package (<https://github.com/singmann/afex>).

²⁰⁴ To control the false positive rate when conducting the post-hoc comparisons, we used
²⁰⁵ Bonferroni correction.

²⁰⁶ ***Bayesian hierarchical generalized linear model (GLM).***

²⁰⁷ The classic NHST approach may ignore the uncertainty in estimate of the parameters
²⁰⁸ for SDT (Rouder & Lu, 2005), and using mean RT assumes normal distribution of RT
²⁰⁹ data, which is always not true because RTs distribution is skewed (Rousselet & Wilcox,
²¹⁰ 2019). To better estimate the uncertainty and use a more appropriate model, we also tried
²¹¹ Bayesian hierarchical generalized linear model to analyze each experiment's accuracy and
²¹² RTs data. We used BRMs (Bürkner, 2017) to build the model, which used Stan (Carpenter
²¹³ et al., 2017) to estimate the posterior.

²¹⁴ In the GLM model, we assume that the accuracy of each trial is Bernoulli distributed
²¹⁵ (binomial with 1 trial), with probability p_i that $y_i = 1$.

$$y_i \sim Bernoulli(p_i)$$

- 216 In the perceptual matching task, the probability p_i can then be modeled as a function of
217 the trial type:

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

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

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

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

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

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

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

235 The subjective-specific intercepts ($\beta_0 = -zFAR$) and slopes ($\beta_1 = d'$) are describe

236 by multivariate normal with means and a covariance matrix for the parameters.

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

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

238 ([https://lindeloev.github.io/shiny-rt/#34_\(shifted\)_log-normal](https://lindeloev.github.io/shiny-rt/#34_(shifted)_log-normal)). This distribution has

239 two parameters: μ, σ . μ is the mean of the logNormal distribution, and σ is the disperse of

240 the distribution. The log normal distribution can be extended to shifted log normal

241 distribution, with one more parameter: shift, which is the earliest possible response.

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

242 Shifted log-normal distribution:

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

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

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

$$\sigma_j \sim Cauchy()$$

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

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

246 used HDDM (Wiecki, Sofer, & Frank, 2013) to model our RTs and accuracy data. We used

247 the prior implemented in HDDM, that is, informative priors that constrains parameter

248 estimates to be in the range of plausible values based on past literature (Matzke &

249 Wagenmakers, 2009). As reported in Hu et al. (2020), we used the response code approach,

match response were coded as 1 and nonmatch responses were coded as 0. To fully explore all parameters, we allow all four parameters of DDM free to vary. We then extracted the estimation of all the four parameters for each participants for the correlation analyses. However, because the starting point is only related to response (match vs. non-match) but not the valence of the stimuli, we didn't included it in correlation analysis.

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

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

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

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

M_1 is the mean of the first condition, sd_1 is the standard deviation of the first condition, while M_2 is the mean of the second condition, sd_2 is the standard deviation of the second condition. r is the correlation coefficient between data from first and second

²⁷¹ condition. n is the number of data point (in our case the number of participants included
²⁷² in our research).

²⁷³ The effect size from each experiment were then synthesized by random effect model
²⁷⁴ using `metafor` (Viechtbauer, 2010). Note that to avoid the cases that some participants
²⁷⁵ participated more than one experiments, we inspected the all available information of
²⁷⁶ participants and only included participants' results from their first participation. As
²⁷⁷ mentioned above, 24 participants were intentionally recruited to participate both exp 1a
²⁷⁸ and exp 2, we only included their results from experiment 1a in the meta-analysis.

²⁷⁹ We also estimated the synthesized effect size using Bayesian hierarchical model,
²⁸⁰ which extended the two-level hierarchical model in each experiment into three-level model,
²⁸¹ which experiment as an additional level. For SDT, we can use a nested hierarchical model
²⁸² to model all the experiment with similar design:

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

²⁸³ where

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

²⁸⁴ The outcomes y_{ijk} are 0 if participant j in experiment k responded “nonmatch” on trial i ,
²⁸⁵ 1 if they responded “match”.

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

²⁸⁶ and the experiment level parameter μ_{0k} and μ_{1k} is from a higher order
²⁸⁷ distribution:

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

²⁸⁸ in which μ_0 and μ_1 means the population level parameter.

289 This model can be easily expand to three-level model in which participants and

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

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

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

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

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

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

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

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

296 ***Valence effect.***

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

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

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

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

304 In the third part, we focused on experiment 4a and 4b, which were designed to
305 examine the implicit effect of the interaction between moral valence and self-referential
306 processing. We are interested in one particular question: will self-referential and morally
307 positive valence had a mutual facilitation effect. That is, when moral valence (experiment
308 4a) or self-referential (experiment 4a) was presented as task-irrelevant stimuli, whether
309 they would facilitate self-referential or valence effect on perceptual decision-making. For
310 experiment 4a, we reported the comparisons between different valence conditions under the
311 self-referential task and other-referential task. For experiment 4b, we first calculated the
312 effect of valence for both self- and other-referential conditions and then compared the effect
313 size of these three contrast from self-referential condition and from other-referential
314 condition. Note that the results were also analyzed in a standard repeated measure
315 ANOVA (see supplementary materials).

316 ***Specificity of the valence effect.***

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

321 ***Behavior-Questionnaire correlation.***

322 Finally, we explored correlation between results from behavioral results and
323 self-reported measures.

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

329 For the behavioral task part, we used three parameters from drift diffusion model:

330 drift rate (v), boundary separation (a), and non decision-making time (t), because these
331 parameters has relative clear psychological meaning. We used the mean of parameter
332 posterior distribution as the estimate of each parameter for each participants in the
333 correlation analysis.

334 Based on results form the experiment, we reason that the correlation between

335 behavioral result in self-referential will appear in the data without mentioning the
336 self/other (exp 3a, 3b, 6a, 7a, 7b). To this end, we first calculated the correlation between
337 behavioral indicators and questionnaires for self-referential and other-referential separately.
338 Given the small sample size of the data ($N =$), we used a relative liberal threshold for
339 these exploration ($\alpha = 0.1$).

340 Then we confirmed the significant results from the data without self- and

341 other-referential (exp1a, 1b, 1c, 2, 5, 6a). In this confirmatory analysis, we used $\alpha =$
342 0.05 and used bootstrap by BootES package (Kirby & Gerlanc, 2013) to estimate the
343 correlation. To avoid false positive, we further determined the threshold for significant by
344 permutation. More specifically, for each pairs that initially with $p < .05$, we randomly
345 shuffle the participants data of each score and calculated the correlation between the
346 shuffled vectors. After repeating this procedure for 5000 times, we choose arrange these
347 5000 correlation coefficients and use the 95% percentile number as our threshold.

348

Part 1: Moral valence effect

349 In this part, we report five experiments that aimed at testing whether the instantly

350 acquired association between shapes and good person would be prioritized in perceptual
351 decision-making.

352 **Experiment 1a**

353 **Methods.**

354 ***Participants.***

355 57 college students (38 female, age = 20.75 ± 2.54 years) participated. 39 of them
356 were recruited from Tsinghua University community in 2014; 18 were recruited from
357 Wenzhou University in 2017. All participants were right-handed except one, and all had
358 normal or corrected-to-normal vision. Informed consent was obtained from all participants
359 prior to the experiment according to procedures approved by the local ethics committees. 6
360 participant's data were excluded from analysis because nearly random level of accuracy,
361 leaving 51 participants (34 female, age = 20.72 ± 2.44 years).

362 ***Stimuli and Tasks.***

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

366 ***Procedure.***

367 This experiment had two phases. First, there was a brief learning stage. Participants
368 were asked to learn the relationship between geometric shapes (triangle, square, and circle)
369 and different person (bad person, a good person, or a neutral person). For example, a
370 participant was told, "bad person is a circle; good person is a triangle; and a neutral person
371 is represented by a square." After participant remember the associations (usually in a few
372 minutes), participants started a practicing phase of matching task which has the exact task
373 as in the experimental task. In the experimental task, participants judged whether
374 shape-label pairs, which were subsequently presented, were correct. Each trial started with
375 the presentation of a central fixation cross for 500 ms. Subsequently, a pairing of a shape
376 and label (good person, bad person, and neutral person) was presented for 100 ms. The

377 pair presented could confirm to the verbal instruction for each pairing given in the training
378 stage, or it could be a recombination of a shape with a different label, with the shape–label
379 pairings being generated at random. The next frame showed a blank for 1100ms.
380 Participants were expected to judge whether the shape was correctly assigned to the person
381 by pressing one of the two response buttons as quickly and accurately as possible within
382 this timeframe (to encourage immediate responding). Feedback (correct or incorrect) was
383 given on the screen for 500 ms at the end of each trial, if no response detected, “too slow”
384 was presented to remind participants to accelerate. Participants were informed of their
385 overall accuracy at the end of each block. The practice phase finished and the experimental
386 task began after the overall performance of accuracy during practice phase achieved 60%.
387 For participants from the Tsinghua community, they completed 6 experimental blocks of 60
388 trials. Thus, there were 60 trials in each condition (bad-person match, bad-person
389 nonmatch, good-person match, good-person nonmatch, neutral-person match, and
390 neutral-person nonmatch). For the participants from Wenzhou University, they finished 6
391 blocks of 120 trials, therefore, 120 trials for each condition.

392 ***Data analysis.***

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

396 **Results.**

397 ***Classic NHST.***

398 *d prime.*

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

401 We found the effect of Valence ($F(1.96, 97.84) = 6.19$, $MSE = 0.27$, $p = .003$,
402 $\hat{\eta}_G^2 = .020$). The post-hoc comparison with multiple comparison correction revealed that

403 the shapes associated with Good-person (2.11, SE = 0.14) has greater d prime than shapes
 404 associated with Bad-person (1.75, SE = 0.14), $t(50) = 3.304$, $p = 0.0049$. The Good-person
 405 condition was also greater than the Neutral-person condition (1.95, SE = 0.16), but didn't
 406 reach statistical significant, $t(50) = 1.54$, $p = 0.28$. Neither the Neutral-person condition is
 407 significantly greater than the Bad-person condition, $t(50) = 2.109$, $p = .098$.

408 *Reaction times.*

409 We conducted 2 (Matchness: match v. nonmatch) by 3 (Valence: good, neutral, bad)
 410 repeated measure ANOVA. We found the main effect of Matchness ($F(1, 50) = 232.39$,
 411 $MSE = 948.92$, $p < .001$, $\hat{\eta}_G^2 = .104$), main effect of valence ($F(1.87, 93.31) = 9.62$,
 412 $MSE = 1,673.86$, $p < .001$, $\hat{\eta}_G^2 = .016$), and interaction between Matchness and Valence
 413 ($F(1.73, 86.65) = 8.52$, $MSE = 1,441.75$, $p = .001$, $\hat{\eta}_G^2 = .011$).

414 We then carried out two separate ANOVA for Match and Mismatched trials. For
 415 matched trials, we found the effect of valence . We further examined the effect of valence
 416 for both self and other for matched trials. We found that shapes associated with Good
 417 Person (684 ms, SE = 11.5) responded faster than Neutral (709 ms, SE = 11.5), $t(50) =$
 418 -2.265, $p = 0.0702$ and Bad Person (728 ms, SE = 11.7), $t(50) = -4.41$, $p = 0.0002$), and
 419 the Neutral condition was faster than the Bad condition, $t(50) = -2.495$, $p = 0.0415$). For
 420 non-matched trials, there was no significant effect of Valence ()�.

421 ***Bayesian hierarchical GLM.***

422 *d prime.*

423 We fitted a Bayesian hierarchical GLM for signal detection theory approach. The
 424 results showed that when the shapes were tagged with labels with different moral valence,
 425 the sensitivity (d') and criteria (c) were both influence. For the d' , we found that the
 426 shapes tagged with morally good person (2.46, 95% CI[2.21 2.72]) is greater than shapes
 427 tagged with moral bad (2.07, 95% CI[1.83 2.32]), $P_{PosteriorComparison} = 1$. Shape tagged
 428 with morally good person is also greater than shapes tagged with neutral person (2.23,

429 95% CI[1.95 2.49]), $P_{PosteriorComparison} = 0.97$. Also, the shapes tagged with neutral
 430 person is greater than shapes tagged with morally bad person, $P_{PosteriorComparison} = 0.92$.

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

436 *Reaction times.*

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

445 **HDDM.**

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

453 **Experiment 1b**

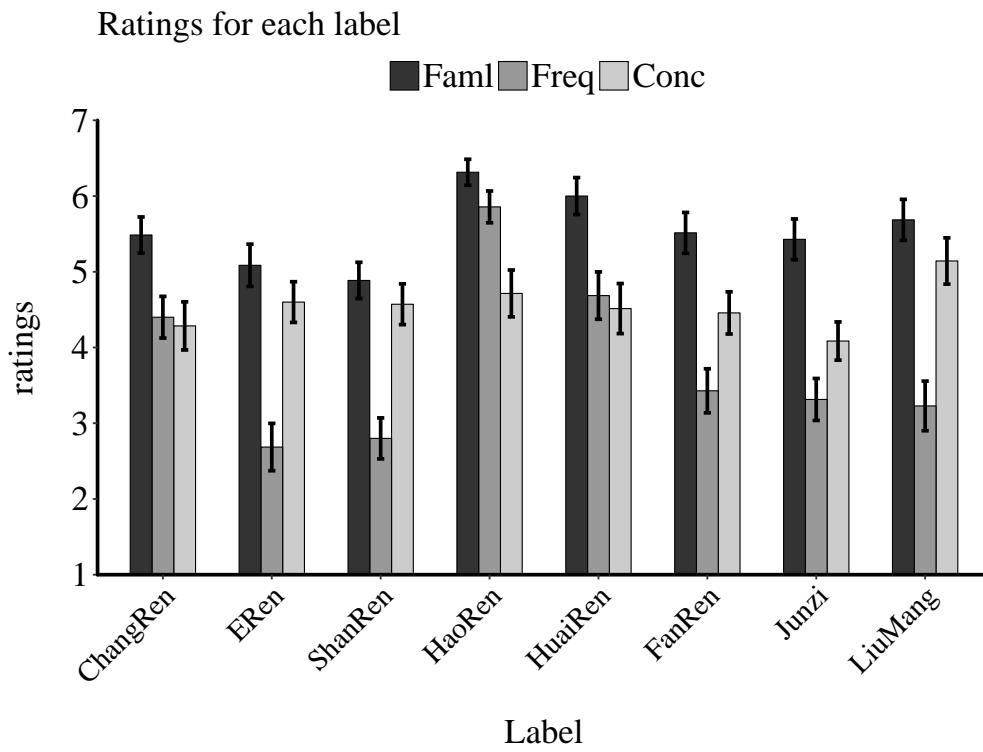
454 In this study, we aimed at excluding the potential confounding factor of the
455 familiarity of words we used in experiment 1a, by matching the familiarity of the words.

456 **Method.**

457 ***Participants.***

458 72 college students (49 female, age = 20.17 ± 2.08 years) participated. 39 of them
459 were recruited from Tsinghua University community in 2014; 33 were recruited from
460 Wenzhou University in 2017. All participants were right-handed except one, and all had
461 normal or corrected-to-normal vision. Informed consent was obtained from all participants
462 prior to the experiment according to procedures approved by the local ethics committees.
463 20 participant's data were excluded from analysis because nearly random level of accuracy,
464 leaving 52 participants (36 female, age = 20.25 ± 2.31 years).

465 **Stimuli and Tasks.** Three geometric shapes (triangle, square, and circle, with 3.7°
466 $\times 3.7^\circ$ of visual angle) were presented above a white fixation cross subtending $0.8^\circ \times 0.8^\circ$
467 of visual angle at the center of the screen. The three shapes were randomly assigned to
468 three labels with different moral valence: a morally bad person (" ", ERen), a morally
469 good person (" ", ShanRen) or a morally neutral person (" ", ChangRen). The order of
470 the associations between shapes and labels was counterbalanced across participants. Three
471 labels used in this experiment is selected based on the rating results from an independent
472 survey, in which participants rated the familiarity, frequency, and concreteness of eight
473 different words online. Of the eight words, three of them are morally positive (HaoRen,
474 ShanRen, Junzi), two of them are morally neutral (ChangRen, FanRen), and three of them
475 are morally negative (HuaiRen, ERen, LiuMang). An independent sample consist of 35
476 participants (22 females, age 20.6 ± 3.11) were recruited to rate these words. Based on the
477 ratings (see supplementary materials Figure S1), we selected ShanRen, ChangRen, and
478 ERen to represent morally positive, neutral, and negative person.



479

Procedure.

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

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

484 **Results.**

485 **NHST.**

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

487 d prime.

488 Repeated measures ANOVA revealed main effect of valence, $F(1.83, 93.20) = 14.98$,
 489 $MSE = 0.18$, $p < .001$, $\hat{\eta}_G^2 = .053$. Paired t test showed that the Good-Person condition
 490 (1.87 ± 0.102) was with greater d prime than Neutral condition (1.44 ± 0.101 , $t(51) =$
 491 5.945 , $p < 0.001$). We also found that the Bad-Person condition (1.67 ± 0.11) has also
 492 greater d prime than neutral condition , $t(51) = 3.132$, $p = 0.008$). There Good-person

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

494 *Reaction times.*

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

496 $MSE = 960.92, p < .001, \hat{\eta}_G^2 = .031$) and then analyzed the matched trials and

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

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

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

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

501 $t(51) = -5.724, p < 0.0001$). While there was no significant differences between Neutral and

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

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

504 **BGLM.**

505 *Signal detection theory analysis of accuracy.*

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

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

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

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

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

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

512 $P_{PosteriorComparison} = 0.97$. Also, the shapes tagged with neutral person is greater than

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

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

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

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

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

518 evidence for the difference between good and bad conditions.

519 *Reaction time.*

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

528 **HDDM.**

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

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

539 **Experiment 1c**

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

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

544 **Method.**

545 ***Participants.***

546 23 college students (15 female, age = 22.61 ± 2.62 years) participated. All of them
547 were recruited from Tsinghua University community in 2014. Informed consent was
548 obtained from all participants prior to the experiment according to procedures approved by
549 the local ethics committees. No participant was excluded because they overall accuracy
550 were above 0.6.

551 ***Stimuli and Tasks.***

552 Three geometric shapes (triangle, square, and circle, with $3.7^\circ \times 3.7^\circ$ of visual angle)
553 were presented above a white fixation cross subtending $0.8^\circ \times 0.8^\circ$ of visual angle at the
554 center of the screen. The three most common names were chosen, which are neutral in
555 moral valence before the manipulation. Three names (Zhang, Wang, Li) were first paired
556 with three paragraphs of behavioral description. Each description includes one sentence of
557 biographic information and four sentences that describing the moral behavioral under that
558 name. To assess the that these three descriptions represented good, neutral, and bad
559 valence, we collected the ratings of three person on six dimensions: morality, likability,
560 trustworthiness, dominance, competence, and aggressiveness, from an independent sample
561 ($n = 34$, 18 female, age = 19.6 ± 2.05). The rating results showed that the person with
562 morally good behavioral description has higher score on morality ($M = 3.59$, $SD = 0.66$)
563 than neutral ($M = 0.88$, $SD = 1.1$), $t(33) = 12.94$, $p < .001$, and bad conditions ($M = -3.4$,
564 $SD = 1.1$), $t(33) = 30.78$, $p < .001$. Neutral condition was also significant higher than bad
565 conditions $t(33) = 13.9$, $p < .001$ (See supplementary materials).

566 ***Procedure.***

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

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

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

580 **Results.** Figure 7 shows d prime and reaction times of experiment 1c. We
 581 conducted same analysis as in Experiment 1a. Our analysis didn't show effect of valence
 582 on d prime, $F(1.93, 42.56) = 0.23$, $MSE = 0.41$, $p = .791$, $\hat{\eta}_G^2 = .005$. Neither the effect of
 583 valence on RT ($F(1.63, 35.81) = 0.22$, $MSE = 2,212.71$, $p = .761$, $\hat{\eta}_G^2 = .001$) or
 584 interaction between valence and matchness on RT ($F(1.79, 39.43) = 1.20$,
 585 $MSE = 1,973.91$, $p = .308$, $\hat{\eta}_G^2 = .005$).

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

587 We fitted a Bayesian hierarchical GLM for SDT. The results showed that when the
 588 shapes were tagged with labels with different moral valence, the sensitivity (d') and criteria
 589 (c) were both influenced. For the d' , we found that the shapes tagged with morally good
 590 person (2.30, 95% CI[1.93 2.70]) is greater than shapes tagged with moral bad (2.11, 95%
 591 CI[1.83 2.42]), $P_{PosteriorComparison} = 0.8$. Shape tagged with morally good person is also
 592 greater than shapes tagged with neutral person (2.16, 95% CI[1.88 2.45]),
 593 $P_{PosteriorComparison} = 0.75$.

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

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

599 ***Reaction time.***

600 We fitted a Bayesian hierarchical GLM for RTs, with a log-normal distribution as the
601 link function. We used the posterior distribution of the regression coefficient to make
602 statistical inferences. As in previous studies, the matched conditions are much faster than
603 the mismatched trials ($P_{PosteriorComparison} = .75$). We focused on matched trials only, and
604 compared different conditions: Good () is not faster than the neutral (),
605 $P_{PosteriorComparison} = .5$, it was faster than the Bad condition (),
606 $P_{PosteriorComparison} = .88$. And the neutral condition is faster than the bad condition,
607 $P_{PosteriorComparison} = .95$. However, the mismatched trials are largely overlapped.

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

615 **Experiment 2: Sequential presenting**

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

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

628 **Method.**

629 ***Participants.***

630 35 participants (17 female, age = 21.66 ± 3.03) were recruited. 24 of them had
631 participated in Experiment 1a (9 male, mean age = 21.9, s.d. = 2.9), and the time gap
632 between these experiment 1a and experiment 2 is at least six weeks. The results of 1
633 participants were excluded from analysis because of less than 60% overall accuracy,
634 remains 34 participants (17 female, age = 21.74 ± 3.04).

635 ***Procedure.***

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

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

649 ***Data analysis.***

650 Data was analyzed as in study 1a.

651 **Results.**

652 ***NHST.***

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

655 d prime.

656 There was evidence for the main effect of valence, $F(1.83, 60.36) = 14.41$,
657 $MSE = 0.23$, $p < .001$, $\hat{\eta}_G^2 = .066$. Paired t test showed that the Good-Person condition
658 (2.79 ± 0.17) was with greater d prime than Netural condition (2.21 ± 0.16 , $t(33) = 4.723$,
659 $p = 0.001$) and Bad-person condition (2.41 ± 0.14), $t(33) = 4.067$, $p = 0.008$). There was
660 no-significant difference between Neutral-person and Bad-person conidition, $t(33) = -1.802$,
661 $p = 0.185$.

662 *Reaction time.*

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

665 We found interaction between Matchness and Valence ($F(1.99, 65.70) = 9.53$,
666 $MSE = 605.36$, $p < .001$, $\hat{\eta}_G^2 = .017$) and then analyzed the matched trials and
667 mismatched trials separately, as in experiment 1a. For matched trials, we found the effect
668 of valence $F(1.99, 65.76) = 10.57$, $MSE = 1,192.65$, $p < .001$, $\hat{\eta}_G^2 = .067$. Post-hoc t -tests
669 revealed that shapes associated with Good Person (548 ± 9.4) were responded faster than
670 Neutral-Person (582 ± 10.9), ($t(33) = -3.95$, $p = 0.0011$) and Bad Person (582 ± 10.2),
671 $t(33) = -3.9$, $p = 0.0013$). While there was no significant differences between Neutral and

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

674 **BGLMM.**

675 *Signal detection theory analysis of accuracy.*

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

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

689 *Reaction times.*

690 We fitted a Bayesian hierarchical GLM for RTs, with a log-normal distribution as the
 691 link function. We used the posterior distribution of the regression coefficient to make
 692 statistical inferences. As in previous studies, the matched conditions are much faster than
 693 the mismatched trials ($P_{PosteriorComparison} = .75$). We focused on matched trials only, and
 694 compared different conditions: Good () is not faster than the neutral (),
 695 $P_{PosteriorComparison} = .5$, it was faster than the Bad condition (),
 696 $P_{PosteriorComparison} = .88$. And the neutral condition is faster than the bad condition,
 697 $P_{PosteriorComparison} = .95$. However, the mismatched trials are largely overlapped.

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

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

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

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

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

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

704 that shapes tagged with bad person had longer non-decision time (see figure

705 @ref(fig:plot-exp1c -HDDM))).

706 Discussion

707 In this experiment, we repeated the results pattern that the positive moral valenced

708 stimuli has an advantage over the neutral or the negative valence association. Moreover,

709 with a cross-task analysis, we did not find evidence that the experiment task interacted

710 with moral valence, suggesting that the effect might not be effect by experiment task.

711 These findings suggested that the facilitation effect of positive moral valence is robust and

712 not affected by task. This robust effect detected by the associative learning is unexpected.

713 Experiment 6a: EEG study 1

714 Experiment 6a was conducted to study the neural correlates of the positive

715 prioritization effect. The behavioral paradigm is same as experiment 2.

716 Method.

717 Participants.

718 24 college students (8 female, age = 22.88 ± 2.79) participated the current study, all

719 of them were from Tsinghua University in 2014. Informed consent was obtained from all

720 participants prior to the experiment according to procedures approved by a local ethics

721 committee. No participant was excluded from behavioral analysis.

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

726 *Stimuli.*

727 Three geometric shapes (triangle, square and circle, each $4.6^\circ \times 4.6^\circ$ of visual angle)
728 were presented at the center of screen for 50 ms after 500ms of fixation ($0.8^\circ \times 0.8^\circ$ of
729 visual angle). The association of the three shapes to bad person (“ , HuaiRen”), good
730 person (“ , HaoRen”) or ordinary person (“ , ChangRen”) was counterbalanced across
731 participants. The words bad person, good person or ordinary person ($3.6^\circ \times 1.6^\circ$) was also
732 displayed at the center fo the screen. Participants had to judge whether the pairings of
733 label and shape matched (e.g., Does the circle represent a bad person?). The experiment
734 was run on a PC using E-prime software (version 2.0). These stimuli were displayed on a
735 22-in CRT monitor (1024×768 at 100Hz). We used backward masking to avoid
736 over-processing of the moral words, in which a scrambled picture were presented for 900 ms
737 after the label. Also, to avoid the ceiling effect on accuracy, shapes were presented on a
738 noisy background based on our pilot studies. The noisy images were made by scrambling a
739 picture of 3/4gray and 1/4 white at resolution of 2×2 pixel.

740 *Procedure.*

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

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

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

752 All the stimuli were presented on a gray background (RGB: 127, 127, 127). E-primed
753 2.0 was used to present stimuli and collect behavioral results. Data were collected and
754 analyzed when accuracy performance in total reached 60%.

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

756 **Results.**

757 **NHST.**

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

760 *d prime.*

761 We conducted repeated measures ANOVA, with moral valence as independent
762 variable. The results revealed the main effect of valence ($F(1.74, 40.05) = 3.76$,
763 $MSE = 0.10$, $p = .037$, $\eta^2_G = .021$). Post-hoc analysis revealed that shapes link with Good
764 person (mean = 3.13, SE = 0.109) is greater than Neutral condition (mean = 2.88, SE =
765 0.14), $t = 2.916$, $df = 24$, $p = 0.02$, p-value adjusted by Tukey method, but the *d* prime
766 between Good and bad (mean = 3.03, SE = 0.142) ($t = 1.512$, $df = 24$, $p = 0.3034$, p-value
767 adjusted by Tukey method), bad and neutral ($t = 1.599$, $df = 24$, $p = 0.2655$, p-value
768 adjusted by Tukey method) were not significant.

769 *Reaction times.*

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

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

773 $MSE = 450.47, p < .001, \hat{\eta}_G^2 = .021$) and then analyzed the matched trials and
 774 mismatched trials separately, as in experiment 2. For matched trials, we found the effect of
 775 valence $F(1.97, 45.25) = 32.37, MSE = 522.42, p < .001, \hat{\eta}_G^2 = .078$. For non-matched
 776 trials, there was no significant effect of Valence ($F(1.77, 40.67) = 0.35, MSE = 242.15,$
 777 $p = .679, \hat{\eta}_G^2 = .000$). Post-hoc t -tests revealed that shapes associated with Good Person
 778 (mean = 550, SE = 13.8) were responded faster than Neutral-Person (501, SE = 14.7),
 779 ($t(24) = -5.171, p = 0.0001$) and Bad Person (523, SE = 16.3), $t(24) = -8.137, p <$
 780 0.0001), and Neutral is faster than Bad-Person condition ($t(32) = -3.282, p = 0.0085$).

781 **BGLM.**

782 *Signal detection theory analysis of accuracy.*

783 *Reaction time.*

784 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et
 785 al., 2020). We estimated separate drift rate (v), non-decision time (T_0), and boundary
 786 separation (a) for each condition. We found that, similar to experiment 2, the shapes
 787 tagged with good person has higher drift rate and higher boundary separation than shapes
 788 tagged with both neutral and bad person, but only for the self-referential condition. Also,
 789 the shapes tagged with neutral person has a higher drift rate than shapes tagged with bad
 790 person, but not for the boundary separation, and this effect also exist only for the
 791 self-referential condition.

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

796

Part 2: interaction between valence and identity

797

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

800

Experiment 3a

801

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

805

Method.

806

Participants.

807

38 college students (15 female, age = 21.92 ± 2.16) participated in experiment 3a.

808

All of them were right-handed, and all had normal or corrected-to-normal vision. Informed
809 consent was obtained from all participants prior to the experiment according to procedures
810 approved by a local ethics committee. One female and one male student did not finish the
811 experiment, and 1 participants' data were excluded from analysis because less than 60%
812 overall accuracy, remains 35 participants (13 female, age = 22.11 ± 2.13).

813

Design.

814

Study 3a combined moral valence with self-relevance, hence the experiment has a $2 \times$
815 3×2 within-subject design. The first variable was self-relevance, include two levels:
816 self-relevance vs. stranger-relevance; the second variable was moral valence, include good,
817 neutral and bad; the third variable was the matching between shape and label: match
818 vs. nonmatch.

819 *Stimuli.*

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

825 *Procedure.*

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

841 *Data Analysis.*

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

845 **Results.**

846 **NHST.**

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

849 *d prime.*

850 There was evidence for the main effect of valence, $F(1.89, 64.37) = 11.09$,
 851 $MSE = 0.23$, $p < .001$, $\hat{\eta}_G^2 = .039$, and main effect of self-relevance, $F(1, 34) = 3.22$,
 852 $MSE = 0.54$, $p = .082$, $\hat{\eta}_G^2 = .015$, as well as the interaction, $F(1.79, 60.79) = 3.39$,
 853 $MSE = 0.43$, $p = .045$, $\hat{\eta}_G^2 = .022$.

854 We then conducted separated ANOVA for self-referential and other-referential trials.
 855 The valence effect was shown for the self-referential conditions, $F(1.65, 56.25) = 13.98$,
 856 $MSE = 0.31$, $p < .001$, $\hat{\eta}_G^2 = .119$. Post-hoc test revealed that the Good-Self condition
 857 (1.97 ± 0.14) was with greater d prime than Neutral condition (1.41 ± 0.12 , $t(34) = 4.505$,
 858 $p = 0.0002$), and Bad-self condition (1.43 ± 0.102), $t(34) = 3.856$, $p = 0.0014$. There was
 859 difference between neutral and bad condition, $t(34) = -0.238$, $p = 0.9694$. However, no
 860 effect of valence was found for the other-referential condition $F(1.98, 67.36) = 0.38$,
 861 $MSE = 0.35$, $p = .681$, $\hat{\eta}_G^2 = .004$.

862 *Reaction time.*

863 We found interaction between Matchness and Valence ($F(1.98, 67.44) = 26.29$,
 864 $MSE = 730.09$, $p < .001$, $\hat{\eta}_G^2 = .025$) and then analyzed the matched trials and nonmatch
 865 trials separately, as in previous experiments.

866 For the match trials, we found that the interaction between identity and valence,
 867 $F(1.72, 58.61) = 3.89$, $MSE = 2,750.19$, $p = .032$, $\hat{\eta}_G^2 = .019$, as well as the main effect of
 868 valence $F(1.98, 67.34) = 35.76$, $MSE = 1,127.25$, $p < .001$, $\hat{\eta}_G^2 = .079$, but not the effect of
 869 identity $F(1, 34) = 0.20$, $MSE = 3,507.14$, $p = .660$, $\hat{\eta}_G^2 = .001$. As for the d prime, we

870 separated analyzed the self-referential and other-referential trials. For the Self-referential
 871 trials, we found the main effect of valence, $F(1.80, 61.09) = 30.39$, $MSE = 1,584.53$,
 872 $p < .001$, $\hat{\eta}_G^2 = .159$; for the other-referential trials, the effect of valence is weaker,
 873 $F(1.86, 63.08) = 2.85$, $MSE = 2,224.30$, $p = .069$, $\hat{\eta}_G^2 = .024$. We then focused on the self
 874 conditions: the good-self condition (713 ± 12) is faster than neutral- (776 ± 11.8), $t(34) =$
 875 -7.396 , $p < .0001$, and bad-self (772 ± 10.1) conditions, $t(34) = -5.66$, $p < .0001$. But
 876 there is not difference between neutral- and bad-self conditions, $t(34) = 0.481$, $p = 0.881$.

877 For the nonmatch trials, we didn't found any strong effect: identity, $F(1, 34) = 3.43$,
 878 $MSE = 660.02$, $p = .073$, $\hat{\eta}_G^2 = .004$, valence $F(1.89, 64.33) = 0.40$, $MSE = 444.10$,
 879 $p = .661$, $\hat{\eta}_G^2 = .001$, or interaction between the two $F(1.94, 66.02) = 2.42$, $MSE = 817.35$,
 880 $p = .099$, $\hat{\eta}_G^2 = .007$.

881 **BGLM.**

882 *Signal detection theory analysis of accuracy.*

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

889 *Reaction time.*

890 In reaction times, we found that same trends in the match trials as in the RT: while
 891 the shapes associated with good self was greater than with good other (log mean diff =
 892 -0.02858 , 95%HPD[-0.070898, 0.0154]), the direction is reversed for neutral and negative
 893 condition. see Figure 17

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

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

901 Experiment 3b

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

907 Participants.

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

Design.

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

Stimuli.

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

Procedure.

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

948 ***Data Analysis.***

949 Same as experiment 3a.

950 **Results.**

951 ***NHST.***

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

954 *d prime.*

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

956 $MSE = 0.33$, $p = .157$, $\hat{\eta}_G^2 = .005$, but we found a main effect of self-relevance,

957 $F(1, 55) = 4.65$, $MSE = 0.89$, $p = .035$, $\hat{\eta}_G^2 = .017$, as well as the interaction,

958 $F(1.90, 104.36) = 5.58$, $MSE = 0.26$, $p = .006$, $\hat{\eta}_G^2 = .011$.

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

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

961 $MSE = 0.30$, $p = .003$, $\hat{\eta}_G^2 = .037$. Post-hoc test revealed that the Good-Self condition

962 (2.15 ± 0.12) was with greater d prime than Neutral condition $(1.83 \pm 0.12$, $t(34) = 3.36$,

963 $p = 0.0031$), and Bad-self condition (1.87 ± 0.12) , $t(34) = 2.955$, $p = 0.01$. There was

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

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

966 $MSE = 0.31$, $p = .539$, $\hat{\eta}_G^2 = .002$.

967 *Reaction time.*

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

969 $MSE = 3, 112.78$, $p < .001$, $\hat{\eta}_G^2 = .006$) and then analyzed the matched trials and

970 nonmatch trials separately, as in previous experiments.

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

972 $F(1.67, 92.11) = 6.14$, $MSE = 6, 472.48$, $p = .005$, $\hat{\eta}_G^2 = .009$, as well as the main effect of

973 valence $F(1.88, 103.65) = 24.25$, $MSE = 5,994.25$, $p < .001$, $\hat{\eta}_G^2 = .038$, but not the effect
 974 of identity $F(1, 55) = 48.49$, $MSE = 25,892.59$, $p < .001$, $\hat{\eta}_G^2 = .153$. As for the d prime,
 975 we separated analyzed the self-referential and other-referential trials. For the
 976 Self-referential trials, we found the main effect of valence, $F(1.66, 91.38) = 23.98$,
 977 $MSE = 6,965.61$, $p < .001$, $\hat{\eta}_G^2 = .100$; for the other-referential trials, the effect of valence
 978 is weaker, $F(1.89, 103.94) = 5.96$, $MSE = 5,589.90$, $p = .004$, $\hat{\eta}_G^2 = .014$. We then focused
 979 on the self conditions: the good-self condition (713 ± 12) is faster than neutral- ($776 \pm$
 980 11.8), $t(34) = -7.396$, $p < .0001$, and bad-self (772 ± 10.1) conditions, $t(34) = -5.66$, $p <$
 981 $.0001$. But there is not difference between neutral- and bad-self conditions, $t(34) = 0.481$, p
 982 $= 0.881$.

983 For the nonmatch trials, we didn't found any strong effect: identity, $F(1, 55) = 10.31$,
 984 $MSE = 24,590.52$, $p = .002$, $\hat{\eta}_G^2 = .035$, valence $F(1.98, 108.63) = 20.57$, $MSE = 2,847.51$,
 985 $p < .001$, $\hat{\eta}_G^2 = .016$, or interaction between the two $F(1.93, 106.25) = 35.51$,
 986 $MSE = 1,939.88$, $p < .001$, $\hat{\eta}_G^2 = .019$.

987 **BGLM.**

988 *Signal detection theory analysis of accuracy.*

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

995 *Reaction time.*

996 In reaction times, we found that same trends in the match trials as in the RT: while
 997 the shapes associated with good self was greater than with good other (log mean diff =
 998 -0.02858 , 95%HPD[-0.070898, 0.0154]), the direction is reversed for neutral and negative

999 condition. see Figure 20

1000 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et
1001 al., 2020). We estimated separate drift rate (v), non-decision time (T_0), and boundary
1002 separation (a) for each condition. We found that, similar to experiment 3a, the shapes
1003 tagged with good person has higher drift rate and higher boundary separation than shapes
1004 tagged with both neutral and bad person, but only for the self-referential condition. Also,
1005 the shapes tagged with neutral person has a higher drift rate than shapes tagged with bad
1006 person, but not for the boundary separation, and this effect also exist only for the
1007 self-referential condition.

1008 Interestingly, we found that in both self-referential and other-referential conditions,
1009 the shapes associated bad valence have higher drift rate and higher boundary separation.
1010 which might suggest that the shape associated with bad stimuli might be prioritized in the
1011 non-match trials (see figure 21)).

1012 **Experiment 6b**

1013 Experiment 6b was conducted to study the neural correlates of the prioritization
1014 effect of positive self, i.e., the neural underlying of the behavioral effect found int
1015 experiment 3a. However, as in experiment 6a, the procedure of this experiment was
1016 modified to adopted to ERP experiment.

1017 **Method.**

1018 ***Participants.***

1019 23 college students (8 female, age = 22.86 ± 2.47) participated the current study, all
1020 of them were recruited from Tsinghua University in 2016. Informed consent was obtained
1021 from all participants prior to the experiment according to procedures approved by a local
1022 ethics committee. For day 1's data, 1 participant was excluded from the current analysis
1023 because of lower than 60% overall accuracy, remaining 22 participants (8 female, age =

1024 22.76 ± 2.49). For day 2's data, one participant dropped out, leaving 22 participants (9
1025 female, age = 23.05 ± 2.46), all of them has overall accuracy higher than 60%.

1026 ***Design.***

1027 The experimental design of this experiment is same as experiment 3: a 2 × 3 × 2
1028 within-subject design with self-relevance (self-relevant vs. other-relevant), moral valence
1029 (good, neutral, and bad) and matchness between shape and label (match vs. mismatch) as
1030 within-subject variables.

1031 ***Stimuli.***

1032 As in experiment 3a, 6 shapes were included (triangle, square, circle, trapezoid,
1033 diamond, regular pentagon), as well as 6 labels (good self, neutral self, bad self, good
1034 person, bad person, neutral person). To match the concreteness of the label, we asked
1035 participant to chosen an unfamiliar name of their own gender to be the stranger.

1036 ***Procedure.***

1037 The procedure was similar to Experiment 2 and 6a. Subjects first learned the
1038 associations between labels and shapes and then completed a shape-label matching task. In
1039 each trial of the matching task, a fixation were first presented for 500 ms, followed by a 50
1040 ms label; then, a scrambled picture presented 900 ms. After the backward mask, the shape
1041 were presented on a noisy background for 50ms. Participant have to response in 1000ms
1042 after the presentation of the shape, and finally, a feedback screen was presented for 500 ms.
1043 The inter-trial interval (ITI) were randomly varied at the range of 1000 ~ 1400 ms.

1044 All the stimuli were presented on a gray background (RGB: 127, 127, 127). E-primed
1045 2.0 was used to present stimuli and collect behavioral results. Data were collected and
1046 analyzed when accuracy performance in total reached 60%.

1047 Because learning 6 associations was more difficult than 3 associations and participant
1048 might have low accuracy (see experiment 3a), the current study had extended to a two-day

1049 paradigm to maximizing the accurate trials that can be used in EEG data. At the first day,
1050 participants learnt the associations and finished 9 blocks of the matching task, each had
1051 120 trials, without EEG recording. That is, each condition has 90 trials.

1052 Participants came back to lab at the second day and finish the same task again, with
1053 EEG recorded. Before the EEG experiment, each participant finished a practice session
1054 again, if their accuracy is equal or higher than 85%, they start the experiment (one
1055 participant used lower threshold 75%). Each participant finished 18 blocks, each has 90
1056 trials. One participant finished additional 6 blocks because of high error rate at the
1057 beginning, another two participant finished addition 3 blocks because of the technique
1058 failure in recording the EEG data. To increase the number of trials that can be used for
1059 EEG data analysis, matched trials has twice number as mismatched trials, therefore, for
1060 matched trials each participants finished 180 trials for each condition, for mismatched
1061 trials, each conditions has 90 trials.

1062 ***Data Analysis.***

1063 Same as experiment 3a.

1064 **Results of Day 1.**

1065 ***NHST.***

1066 Figure 22 shows d prime and reaction times of experiment 3b. Less than 5% correct
1067 trials with less than 200ms reaction times were excluded.

1068 ***d prime.***

1069 There was no evidence for the main effect of valence, $F(1.91, 40.20) = 11.98$,
1070 $MSE = 0.15$, $p < .001$, $\hat{\eta}_G^2 = .040$, but we found a main effect of self-relevance,
1071 $F(1, 21) = 1.21$, $MSE = 0.20$, $p = .284$, $\hat{\eta}_G^2 = .003$, as well as the interaction,
1072 $F(1.28, 26.90) = 12.88$, $MSE = 0.21$, $p = .001$, $\hat{\eta}_G^2 = .041$.

1073 We then conducted separated ANOVA for self-referential and other-referential trials.

1074 The valence effect was shown for the self-referential conditions, $F(1.73, 36.42) = 29.31$,
 1075 $MSE = 0.14$, $p < .001$, $\hat{\eta}_G^2 = .147$. Post-hoc test revealed that the Good-Self condition
 1076 (2.15 ± 0.12) was with greater d prime than Neutral condition (1.83 ± 0.12 , $t(34) = 3.36$,
 1077 $p = 0.0031$), and Bad-self condition (1.87 ± 0.12), $t(34) = 2.955$, $p = 0.01$. There was
 1078 difference between neutral and bad condition, $t(34) = -0.039$, $p = 0.914$. However, no effect
 1079 of valence was found for the other-referential condition $F(1.75, 36.72) = 0.00$, $MSE = 0.18$,
 1080 $p = .999$, $\hat{\eta}_G^2 = .000$.

1081 *Reaction time.*

1082 We found interaction between Matchness and Valence ($F(1.79, 37.63) = 4.07$,
 1083 $MSE = 704.90$, $p = .029$, $\hat{\eta}_G^2 = .003$) and then analyzed the matched trials and nonmatch
 1084 trials separately, as in previous experiments.

1085 For the match trials, we found that the interaction between identity and valence,
 1086 $F(1.72, 36.16) = 4.55$, $MSE = 1,560.90$, $p = .022$, $\hat{\eta}_G^2 = .015$, as well as the main effect of
 1087 valence $F(1.93, 40.55) = 9.83$, $MSE = 1,951.84$, $p < .001$, $\hat{\eta}_G^2 = .044$, but not the effect of
 1088 identity $F(1, 21) = 4.87$, $MSE = 2,032.05$, $p = .039$, $\hat{\eta}_G^2 = .012$. As for the d prime, we
 1089 separated analyzed the self-referential and other-referential trials. For the Self-referential
 1090 trials, we found the main effect of valence, $F(1.92, 40.38) = 14.48$, $MSE = 1,647.20$,
 1091 $p < .001$, $\hat{\eta}_G^2 = .112$; for the other-referential trials, the effect of valence is weaker,
 1092 $F(1.79, 37.50) = 1.04$, $MSE = 1,842.07$, $p = .356$, $\hat{\eta}_G^2 = .008$. We then focused on the self
 1093 conditions: the good-self condition (713 ± 12) is faster than neutral- (776 ± 11.8), $t(34) =$
 1094 -7.396 , $p < .0001$, and bad-self (772 ± 10.1) conditions, $t(34) = -5.66$, $p < .0001$. But
 1095 there is not difference between neutral- and bad-self conditions, $t(34) = 0.481$, $p = 0.881$.

1096 For the nonmatch trials, we didn't found any strong effect: identity, $F(1, 21) = 2.76$,
 1097 $MSE = 1,718.93$, $p = .112$, $\hat{\eta}_G^2 = .006$, valence $F(1.61, 33.77) = 3.81$, $MSE = 1,532.21$,
 1098 $p = .041$, $\hat{\eta}_G^2 = .012$, or interaction between the two $F(1.90, 39.97) = 2.23$, $MSE = 720.80$,
 1099 $p = .123$, $\hat{\eta}_G^2 = .004$.

BGLM.*Signal detection theory analysis of accuracy.*

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

Reaction time.

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

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

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

Part 3: Implicit binding between valence and identity

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

Experiment 4a: Morality as task-irrelevant variable

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

Method.***Participants.***

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

Design.

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

Stimuli.

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

Procedure.

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

₁₁₇₅ 120 trials for 2017 data. Due to procedure error, the data collected in 2015 in Tsinghua
₁₁₇₆ community only have 60 trials for each block, i.e., 30 trials per condition.

₁₁₇₇ As in study 3a, before each task, the instruction showed the meaning of each label to
₁₁₇₈ participants. The self-matching task and other-matching task were randomized between
₁₁₇₉ participants. Each participant finished 6 blocks, each have 120 trials.

₁₁₈₀ ***Data Analysis.***

₁₁₈₁ Same as experiment 3a.

₁₁₈₂ **Results.**

₁₁₈₃ ***NHST.***

₁₁₈₄ Figure 25 shows d prime and reaction times of experiment 3a. Less than 5% correct
₁₁₈₅ trials with less than 200ms reaction times were excluded.

₁₁₈₆ d prime.

₁₁₈₇ There was no evidence for the main effect of valence, $F(1.93, 111.66) = 0.53$,
₁₁₈₈ $MSE = 0.12$, $p = .581$, $\hat{\eta}_G^2 = .000$, but we found a main effect of self-relevance,
₁₁₈₉ $F(1, 58) = 121.04$, $MSE = 0.48$, $p < .001$, $\hat{\eta}_G^2 = .189$, as well as the interaction,
₁₁₉₀ $F(1.99, 115.20) = 4.12$, $MSE = 0.14$, $p = .019$, $\hat{\eta}_G^2 = .004$.

₁₁₉₁ We then conducted separated ANOVA for self-referential and other-referential trials.
₁₁₉₂ The valence effect was shown for the self-referential conditions, $F(1.95, 112.92) = 3.01$,
₁₁₉₃ $MSE = 0.15$, $p = .055$, $\hat{\eta}_G^2 = .008$. Post-hoc test revealed that the Good-Self condition
₁₁₉₄ (2.15 ± 0.12) was with greater d prime than Neutral condition (1.83 ± 0.12 , $t(34) = 3.36$,
₁₁₉₅ $p = 0.0031$), and Bad-self condition (1.87 ± 0.12), $t(34) = 2.955$, $p = 0.01$. There was
₁₁₉₆ difference between neutral and bad condition, $t(34) = -0.039$, $p = 0.914$. However, no effect
₁₁₉₇ of valence was found for the other-referential condition $F(1.98, 114.61) = 1.75$,
₁₁₉₈ $MSE = 0.10$, $p = .179$, $\hat{\eta}_G^2 = .003$.

₁₁₉₉ Reaction time.

1200 We found interaction between Matchness and Valence ($F(1.94, 112.64) = 0.84$,
 1201 $MSE = 465.35$, $p = .432$, $\hat{\eta}_G^2 = .000$) and then analyzed the matched trials and nonmatch
 1202 trials separately, as in previous experiments.

1203 For the match trials, we found that the interaction between identity and valence,
 1204 $F(1.90, 110.18) = 4.41$, $MSE = 465.91$, $p = .016$, $\hat{\eta}_G^2 = .003$, as well as the main effect of
 1205 valence $F(1.98, 114.82) = 0.94$, $MSE = 606.30$, $p = .392$, $\hat{\eta}_G^2 = .001$, but not the effect of
 1206 identity $F(1, 58) = 124.15$, $MSE = 4,037.53$, $p < .001$, $\hat{\eta}_G^2 = .257$. As for the d prime, we
 1207 separated analyzed the self-referential and other-referential trials. For the Self-referential
 1208 trials, we found the main effect of valence, $F(1.97, 114.32) = 6.29$, $MSE = 367.25$,
 1209 $p = .003$, $\hat{\eta}_G^2 = .006$; for the other-referential trials, the effect of valence is weaker,
 1210 $F(1.95, 112.89) = 0.35$, $MSE = 699.50$, $p = .699$, $\hat{\eta}_G^2 = .001$. We then focused on the self
 1211 conditions: the good-self condition (713 ± 12) is faster than neutral- (776 ± 11.8), $t(34) =$
 1212 -7.396 , $p < .0001$, and bad-self (772 ± 10.1) conditions, $t(34) = -5.66$, $p < .0001$. But
 1213 there is not difference between neutral- and bad-self conditions, $t(34) = 0.481$, $p = 0.881$.

1214 For the nonmatch trials, we didn't found any strong effect: identity, $F(1, 58) = 0.16$,
 1215 $MSE = 1,547.37$, $p = .692$, $\hat{\eta}_G^2 = .000$, valence $F(1.96, 113.52) = 0.68$, $MSE = 390.26$,
 1216 $p = .508$, $\hat{\eta}_G^2 = .000$, or interaction between the two $F(1.90, 110.27) = 0.04$,
 1217 $MSE = 585.80$, $p = .953$, $\hat{\eta}_G^2 = .000$.

1218 **BGLM.**

1219 *Signal detection theory analysis of accuracy.*

1220 We found that the d prime is greater when shapes were associated with good self
 1221 condition than with neutral self or bad self, but shapes associated with bad self and neutral
 1222 self didn't show differences. comparing the self vs other under three condition revealed that
 1223 shapes associated with good self is greater than with good other, but with a weak evidence.
 1224 In contrast, for both neutral and bad valence condition, shapes associated with other had
 1225 greater d prime than with self.

1226 *Reaction time.*

1227 In reaction times, we found that same trends in the match trials as in the RT: while
1228 the shapes associated with good self was greater than with good other (log mean diff =
1229 -0.02858, 95%HPD[-0.070898, 0.0154]), the direction is reversed for neutral and negative
1230 condition. see Figure 26

1231 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et
1232 al., 2020). We estimated separate drift rate (v), non-decision time (T_0), and boundary
1233 separation (a) for each condition. We found that the shapes tagged with good person has
1234 higher drift rate and higher boundary separation than shapes tagged with both neutral and
1235 bad person. Also, the shapes tagged with neutral person has a higher drift rate than
1236 shapes tagged with bad person, but not for the boundary separation. Finally, we found
1237 that shapes tagged with bad person had longer non-decision time (see figure 27)).

1238 **Experiment 4b: Morality as task-irrelevant variable**

1239 In study 4b, we changed the role of valence and identity in task. In this experiment,
1240 participants learn the association between moral valence and the made perceptual match
1241 judgments to associations between different moral valence and shapes as in study 1-3.
1242 Different from experiment 1 ~ 3, we made put the labels of “self/other” in the shapes so
1243 that identity served as an task irrelevant variable. As in experiment 4b, we also
1244 hypothesized that the intrinsic binding between morally good and self will enhance the
1245 performance of good self condition, even identity is irrelevant to the task.

1246 **Method.**

1247 **Participants.**

1248 53 participants (39 female, age = 20.57 ± 1.81) participated the current study, 34 of
1249 them were from Tsinghua University in 2015, 19 were from Wenzhou University
1250 participated in 2017. All participants were right-handed, and all had normal or

1251 corrected-to-normal vision. Informed consent was obtained from all participants prior to
1252 the experiment according to procedures approved by a local ethics committee. The data
1253 from 8 participants from Wenzhou site were excluded from analysis because their accuracy
1254 was close to chance (< 0.6). The results for the remaining 45 participants (33 female, age
1255 = 20.78 ± 1.76) were analyzed and reported.

1256 ***Design.***

1257 As in Experiment 3, a $2 \times 3 \times 2$ within-subject design was used. The first variable was
1258 self-relevance (self and stranger associations); the second variable was moral valence (good,
1259 neutral and bad associations); the third variable was the matching between shape and label
1260 (matching vs. non-match for the personal association). However, in this the task,
1261 participants only learn the association between two geometric shapes and two labels (self
1262 and other), i.e., only self-relevance were related to the task. The moral valence
1263 manipulation was achieved by embedding the personal label of the labels in the geometric
1264 shapes, see below. For simplicity, the trials where shapes where paired with self and with a
1265 word of “good person” inside were shorted as good-self condition, similarly, the trials where
1266 shapes paired with the self and with a word of “bad person” inside were shorted as bad-self
1267 condition. Hence, we also have six conditions: good-self, neutral-self, bad-self, good-other,
1268 neutral-other, and bad-other.

1269 ***Stimuli.***

1270 2 shapes were included (circle, square) and each appeared above a central fixation
1271 cross with the personal label appearing below. However, the shapes were not empty but
1272 with a two-Chinese-character word in the middle, the word was one of three labels with
1273 different moral valence: “good person”, “bad person” and “neutral person”. Before the
1274 experiment, participants learned the self/other association, and were informed to only
1275 response to the association between shapes’ configure and the labels below the fixation, but
1276 ignore the words within shapes. Besides the behavioral experiments, participants from

1277 Tsinghua community also finished questionnaires as Experiments 3, and participants from
1278 Wenzhou community finished a series of questionnaire as the other experiment finished in
1279 Wenzhou.

1280 ***Procedure.***

1281 The procedure was similar to Experiment 1. There were 6 blocks of trial, each with
1282 120 trials for 2017 data. Due to procedure error, the data collected in 2015 in Tsinghua
1283 community only have 60 trials for each block, i.e., 30 trials per condition.

1284 As in study 3a, before each task, the instruction showed the meaning of each label to
1285 participants. The self-matching task and other-matching task were randomized between
1286 participants. Each participant finished 6 blocks, each have 120 trials.

1287 ***Data Analysis.***

1288 Same as experiment 3a.

1289 **Results.**

1290 ***NHST.***

1291 Figure 28 shows d prime and reaction times of experiment 3a. Less than 5% correct
1292 trials with less than 200ms reaction times were excluded.

1293 ***d prime.***

1294 There was no evidence for the main effect of valence, $F(1.59, 69.94) = 2.34$,
1295 $MSE = 0.48$, $p = .115$, $\hat{\eta}_G^2 = .010$, but we found a main effect of self-relevance,
1296 $F(1, 44) = 0.00$, $MSE = 0.08$, $p = .994$, $\hat{\eta}_G^2 = .000$, as well as the interaction,
1297 $F(1.96, 86.41) = 0.53$, $MSE = 0.10$, $p = .585$, $\hat{\eta}_G^2 = .001$.

1298 We then conducted separated ANOVA for self-referential and other-referential trials.
1299 The valence effect was shown for the self-referential conditions, $F(1.75, 76.86) = 3.08$,
1300 $MSE = 0.25$, $p = .058$, $\hat{\eta}_G^2 = .017$. Post-hoc test revealed that the Good-Self condition
1301 (2.15 ± 0.12) was with greater d prime than Neutral condition (1.83 ± 0.12 , $t(34) = 3.36$,

¹³⁰² $p = 0.0031$), and Bad-self condition (1.87 ± 0.12), $t(34) = 2.955$, $p = 0.01$. There was
¹³⁰³ difference between neutral and bad condition, $t(34) = -0.039$, $p = 0.914$. However, no effect
¹³⁰⁴ of valence was found for the other-referential condition $F(1.63, 71.50) = 1.07$, $MSE = 0.33$,
¹³⁰⁵ $p = .336$, $\hat{\eta}_G^2 = .006$.

¹³⁰⁶ *Reaction time.*

¹³⁰⁷ We found interaction between Matchness and Valence ($F(1.87, 82.50) = 18.58$,
¹³⁰⁸ $MSE = 1,291.12$, $p < .001$, $\hat{\eta}_G^2 = .023$) and then analyzed the matched trials and
¹³⁰⁹ nonmatch trials separately, as in previous experiments.

¹³¹⁰ For the match trials, we found that the interaction between identity and valence,
¹³¹¹ $F(1.86, 81.84) = 5.22$, $MSE = 308.30$, $p = .009$, $\hat{\eta}_G^2 = .003$, as well as the main effect of
¹³¹² valence $F(1.80, 79.37) = 11.04$, $MSE = 2,937.54$, $p < .001$, $\hat{\eta}_G^2 = .059$, but not the effect of
¹³¹³ identity $F(1, 44) = 0.23$, $MSE = 263.26$, $p = .632$, $\hat{\eta}_G^2 = .000$. As for the d prime, we
¹³¹⁴ separated analyzed the self-referential and other-referential trials. For the Self-referential
¹³¹⁵ trials, we found the main effect of valence, $F(1.74, 76.48) = 13.69$, $MSE = 1,732.08$,
¹³¹⁶ $p < .001$, $\hat{\eta}_G^2 = .079$; for the other-referential trials, the effect of valence is weaker,
¹³¹⁷ $F(1.87, 82.44) = 7.09$, $MSE = 1,527.43$, $p = .002$, $\hat{\eta}_G^2 = .043$. We then focused on the self
¹³¹⁸ conditions: the good-self condition (713 ± 12) is faster than neutral- (776 ± 11.8), $t(34) =$
¹³¹⁹ -7.396 , $p < .0001$, and bad-self (772 ± 10.1) conditions, $t(34) = -5.66$, $p < .0001$. But
¹³²⁰ there is not difference between neutral- and bad-self conditions, $t(34) = 0.481$, $p = 0.881$.

¹³²¹ For the nonmatch trials, we didn't found any strong effect: identity, $F(1, 44) = 1.96$,
¹³²² $MSE = 319.47$, $p = .169$, $\hat{\eta}_G^2 = .001$, valence $F(1.69, 74.54) = 6.59$, $MSE = 886.19$,
¹³²³ $p = .004$, $\hat{\eta}_G^2 = .010$, or interaction between the two $F(1.88, 82.57) = 0.31$, $MSE = 316.96$,
¹³²⁴ $p = .718$, $\hat{\eta}_G^2 = .000$.

¹³²⁵ **BGLM.**

¹³²⁶ *Signal detection theory analysis of accuracy.*

¹³²⁷ We found that the d prime is greater when shapes were associated with good self

1328 condition than with neutral self or bad self, but shapes associated with bad self and neutral
1329 self didn't show differences. comparing the self vs other under three condition revealed that
1330 shapes associated with good self is greater than with good other, but with a weak evidence.
1331 In contrast, for both neutral and bad valence condition, shapes associated with other had
1332 greater d prime than with self.

1333 *Reaction time.*

1334 In reaction times, we found that same trends in the match trials as in the RT: while
1335 the shapes associated with good self was greater than with good other (log mean diff =
1336 -0.02858, 95%HPD[-0.070898, 0.0154]), the direction is reversed for neutral and negative
1337 condition. see Figure 29

1338 **HDDM.** We fitted our data with HDDM, using the response-coding (also see Hu et
1339 al., 2020). We estimated separate drift rate (v), non-decision time (T_0), and boundary
1340 separation (a) for each condition. We found that the shapes tagged with good person has
1341 higher drift rate and higher boundary separation than shapes tagged with both neutral and
1342 bad person. Also, the shapes tagged with neutral person has a higher drift rate than
1343 shapes tagged with bad person, but not for the boundary separation. Finally, we found
1344 that shapes tagged with bad person had longer non-decision time (see figure 30)).

1345

Results

1346 **Effect of moral valence**

1347 In this part, we synthesized results from experiment 1a, 1b, 1c, 2, 5 and 6a. Data
1348 from 192 participants were included in these analyses. We found differences between
1349 positive and negative conditions on RT was Cohen's $d = -0.58 \pm 0.06$, 95% CI [-0.70 -0.47];
1350 on d' was Cohen's $d = 0.24 \pm 0.05$, 95% CI [0.15 0.34]. The effect was also observed
1351 between positive and neutral condition, RT: Cohen's $d = -0.44 \pm 0.10$, 95% CI [-0.63
1352 -0.25]; d' : Cohen's $d = 0.31 \pm 0.07$, 95% CI [0.16 0.45]. And the difference between neutral

₁₃₅₃ and bad conditions are not significant, RT: Cohen's $d = 0.15 \pm 0.07$, 95% CI [0.00 0.30];
₁₃₅₄ d' : Cohen's $d = 0.07 \pm 0.07$, 95% CI [-0.08 0.21]. See Figure 31 left panel.

₁₃₅₅ **Interaction between valence and self-reference**

₁₃₅₆ In this part, we combined the experiments that explicitly manipulated the
₁₃₅₇ self-reference and valence, which includes 3a, 3b, 6b, 7a, and 7b. For the positive versus
₁₃₅₈ negative contrast, data were from five experiments with 178 participants; for positive
₁₃₅₉ versus neutral and neutral versus negative contrasts, data were from three experiments ((
₁₃₆₀ 3a, 3b, and 6b) with 108 participants.

₁₃₆₁ In most of these experiments, the interaction between self-reference and valence was
₁₃₆₂ significant (see results of each experiment in supplementary materials). In the
₁₃₆₃ mini-meta-analysis, we analyzed the valence effect for self-referential condition and
₁₃₆₄ other-referential condition separately.

₁₃₆₅ For the self-referential condition, we found the same pattern as in the first part of
₁₃₆₆ results. That is we found significant differences between positive and neutral as well as
₁₃₆₇ positive and negative, but not neutral and negative. The effect size of RT between positive
₁₃₆₈ and negative is Cohen's $d = -0.89 \pm 0.12$, 95% CI [-1.11 -0.66]; on d' was Cohen's $d = 0.61$
₁₃₆₉ ± 0.09 , 95% CI [0.44 0.78]. The effect was also observed between positive and neutral
₁₃₇₀ condition, RT: Cohen's $d = -0.76 \pm 0.13$, 95% CI [-1.01 -0.50]; d' : Cohen's $d = 0.69 \pm$
₁₃₇₁ 0.14, 95% CI [0.42 0.96]. And the difference between neutral and bad conditions are not
₁₃₇₂ significant, RT: Cohen's $d = 0.03 \pm 0.13$, 95% CI [-0.22 0.29]; d' : Cohen's $d = 0.08 \pm 0.08$,
₁₃₇₃ 95% CI [-0.07 0.24]. See Figure 31 the middle panel.

₁₃₇₄ For the other-referential condition, we found that only the difference between positive
₁₃₇₅ and negative on RT was significant, all the other conditions were not. The effect size of RT
₁₃₇₆ between positive and negative is Cohen's $d = -0.28 \pm 0.05$, 95% CI [-0.38 -0.17]; on d' was
₁₃₇₇ Cohen's $d = -0.02 \pm 0.08$, 95% CI [-0.17 0.13]. The effect was not observed between

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

1382 Generalizability of the valence effect

1383 In this part, we reported the results from experiment 4 in which either moral valence
1384 or self-reference were manipulated as task-irrelevant stimuli.

1385 For experiment 4a, when self-reference was the target and moral valence was
1386 task-irrelevant, we found that only under the implicit self-referential condition, i.e., when
1387 the moral words were presented as task irrelevant stimuli, there was the main effect of
1388 valence and interaction between valence and reference for both d prime and RT (See
1389 supplementary results for the detailed statistics). For d prime, we found good-self
1390 condition (2.55 ± 0.86) had higher d prime than bad-self condition (2.38 ± 0.80); good self
1391 condition was also higher than neutral self (2.45 ± 0.78) but there was not statistically
1392 significant, while the neutral-self condition was higher than bad self condition and not
1393 significant neither. For reaction times, good-self condition (654.26 ± 67.09) were faster
1394 relative to bad-self condition (665.64 ± 64.59), and over neutral-self condition ($664.26 \pm$
1395 64.71). The difference between neutral-self and bad-self conditions were not significant.
1396 However, for the other-referential condition, there was no significant differences between
1397 different valence conditions. See Figure 32.

1398 For experiment 4b, when valence was the target and the identity was task-irrelevant,
1399 we found a strong valence effect (see supplementary results and Figure 33, Figure 34).

1400 In this experiment, the advantage of good-self condition can only be disentangled by
1401 comparing the self-referential and other-referential conditions. Therefore, we calculated the
1402 differences between the valence effect under self-referential and other referential conditions

1403 and used the weighted variance as the variance of this differences. We found this
1404 modulation effect on RT. The valence effect of RT was stronger in self-referential than
1405 other-referential for the Good vs. Neutral condition (-0.33 ± 0.01), and to a less extent the
1406 Good vs. Bad condition (-0.17 ± 0.01). While the size of the other effect's CI included
1407 zero, suggesting those effects didn't differ from zero. See Figure 35.

1408 **Specificity of valence effect**

1409 In this part, we analyzed the results from experiment 5, which included positive,
1410 neutral, and negative valence from four different domains: morality, emotion, aesthetics of
1411 human, and aesthetics of scene. We found interaction between valence and domain for both
1412 d prime and RT (match trials). A common pattern appeared in all four domains: each
1413 domain showed a binary results instead of gradient on both d prime and RT. For morality,
1414 aesthetics of human, and aesthetics of scene, there was a positivity effect where the positive
1415 conditions had advantages over both neutral (greater d prime and faster RT), while neutral
1416 and negative conditions didn't differ from each other. But for the emotional stimuli, there
1417 was a reversed negativity effect: positive and neutral conditions were not significantly
1418 different from each other but both had advantage over negative conditions. See
1419 supplementary materials for detailed statistics. Also note that the effect size in moral
1420 domain is smaller than the aesthetic domains (beauty of people and beauty of scene). See
1421 Figure 36.

1422 **Self-reported personal distance**

1423 See Figure 37.

1424 **Correlation analyses**

1425 The reliability of questionnaires can be found in (Liu et al., 2020). We calculated the
1426 correlation between the data from behavioral task and the questionnaire data. First, we
1427 calculated the score for each scale based on their structure and factor loading, instead of
1428 sum score (McNeish & Wolf, 2020). Then, we used SEM to estimate the correlation
1429 because it can include measurement model and statistical model in a unified framework.

1430 To make sure that what we found were not false positive, we used two method to
1431 ensure the robustness of our analysis. first, we split the data into two half: the data with
1432 self and without, then, we used the conditional random forest to find the robust correlation
1433 in the exploratory data (with self reference) that can be replicated in the confirmatory data
1434 (without the self reference). The robust correlation were then analyzed using SEM

1435 Instead of use the exploratory correlation analysis, we used a more principled way to
1436 explore the correlation between parameter of HDDM (v , t , and a) and scale scores and
1437 person distance.

1438 We didn't find the correlation between scale scores and the parameters of HDDM,
1439 but found weak correlation between personal distance and the parameter estimated from
1440 Good and neutral conditions.

1441 First, boundary separation (a) of moral good condition was correlated with both
1442 Self-Bad distance ($r = 0.198$, 95% CI [], $p = 0.0063$) and Neutral-Bad distance
1443 ($r = 0.1571$, 95% CI [], $p = 0.031$). At the same time, the non-decision time is negatively
1444 correlated with Self-Bad distance ($r = 0.169$, 95% CI [], $p = 0.0197$). See Figure 38.

1445 Second, we found the boundary separation of neutral condition is positively
1446 correlated with the personal distance between self and good distance ($r = 0.189$, 95% CI [],
1447 $p = 0.036$), but negatively correlated with self-neutral distance($r = -0.183$, 95% CI [],
1448 $p = 0.042$). Also, the drift rate of the neutral condition is positively correlated with the
1449 Self-Bad distance ($r = 0.177$, 95% CI [], $p = 0.048$).a. See figure 39

1450 We also explored the correlation between behavioral data and questionnaire scores
1451 separately for experiments with and without self-referential, however, the sample size is
1452 very low for some conditions.

1453 **Discussion**

1454 **References**

- 1455 Anderson, E., Siegel, E. H., Bliss-Moreau, E., & Barrett, L. F. (2011). The visual impact
1456 of gossip. *Science*, 332(6036), 1446–1448. <https://doi.org/10.1126/science.1201574>
- 1457 Brady, W. J., Gantman, A. P., & Van Bavel, J. J. (2020). Attentional capture helps
1458 explain why moral and emotional content go viral. *Journal of Experimental
1459 Psychology: General*, 149(4), 746–756. <https://doi.org/10.1037/xge0000673>
- 1460 Brainard, D. H. (1997). The psychophysics toolbox. *Spatial Vision*, 10(4), 433–436.
1461 Journal Article.
- 1462 Bürkner, P.-C. (2017). Brms: An r package for bayesian multilevel models using stan.
1463 *Journal of Statistical Software; Vol 1, Issue 1 (2017)*. Journal Article. Retrieved
1464 from
1465 <https://www.jstatsoft.org/v080/i01> <http://dx.doi.org/10.18637/jss.v080.i01>
- 1466 Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M., ...
1467 Riddell, A. (2017). Stan: A probabilistic programming language. *Journal of
1468 Statistical Software*, 76(1). Journal Article. <https://doi.org/10.18637/jss.v076.i01>
- 1469 Cooper, H., Hedges, L. V., & Valentine, J. C. (2009). *The handbook of research synthesis
1470 and meta-analysis* (2nd ed.). Book, New York: Sage.
- 1471 DeCarlo, L. T. (1998). Signal detection theory and generalized linear models. *Psychological
1472 Methods*, 3(2), 186–205. Journal Article. <https://doi.org/10.1037/1082-989X.3.2.186>

- 1473 Faul, F., Erdfelder, E., Buchner, A., & Lang, A.-G. (2009). Statistical power analyses using
1474 g*Power 3.1: Tests for correlation and regression analyses. *Behavior Research
1475 Methods*, 41(4), 1149–1160. <https://doi.org/10.3758/BRM.41.4.1149>
- 1476 Gantman, A. P., & Van Bavel, J. J. (2014). The moral pop-out effect: Enhanced
1477 perceptual awareness of morally relevant stimuli. *Cognition*, 132(1), 22–29.
1478 <https://doi.org/10.1016/j.cognition.2014.02.007>
- 1479 Goh, J. X., Hall, J. A., & Rosenthal, R. (2016). Mini meta-analysis of your own studies:
1480 Some arguments on why and a primer on how. *Social and Personality Psychology
1481 Compass*, 10(10), 535–549. Journal Article. <https://doi.org/10.1111/spc3.12267>
- 1482 Hu, C.-P., Lan, Y., Macrae, C. N., & Sui, J. (2020). Good me bad me: Does valence
1483 influence self-prioritization during perceptual decision-making? *Collabra:
1484 Psychology*, 6(1), 20. Journal Article. <https://doi.org/10.1525/collabra.301>
- 1485 Juechems, K., & Summerfield, C. (2019). Where does value come from? *Trends in
1486 Cognitive Sciences*, 23(10), 836–850. <https://doi.org/10.1016/j.tics.2019.07.012>
- 1487 Kirby, K. N., & Gerlanc, D. (2013). BootES: An r package for bootstrap confidence
1488 intervals on effect sizes. *Behavior Research Methods*, 45(4), 905–927.
1489 <https://doi.org/10.3758/s13428-013-0330-5>
- 1490 Liu, Q., Wang, F., Yan, W., Peng, K., Sui, J., & Hu, C.-P. (2020). Questionnaire data from
1491 the revision of a chinese version of free will and determinism plus scale. *Journal of
1492 Open Psychology Data*, 8(1), 1. Journal Article. <https://doi.org/10.5334/jopd.49/>
- 1493 Matzke, D., & Wagenmakers, E.-J. (2009). Psychological interpretation of the ex-gaussian
1494 and shifted wald parameters: A diffusion model analysis. *Psychonomic Bulletin &
1495 Review*, 16(5), 798–817. <https://doi.org/10.3758/PBR.16.5.798>
- 1496 McNeish, D., & Wolf, M. G. (2020). Thinking twice about sum scores. *Behavior Research
1497 Methods*. <https://doi.org/10.3758/s13428-020-01398-0>

- 1498 Pelli, D. G. (1997). The videotoolbox software for visual psychophysics: Transforming
1499 numbers into movies. *Spatial Vision*, 10(4), 437–442. Journal Article.
- 1500 Rouder, J. N., & Lu, J. (2005). An introduction to bayesian hierarchical models with an
1501 application in the theory of signal detection. *Psychonomic Bulletin & Review*,
1502 12(4), 573–604. Journal Article. <https://doi.org/10.3758/bf03196750>
- 1503 Rousselet, G. A., & Wilcox, R. R. (2019). Reaction times and other skewed distributions:
1504 Problems with the mean and the median. *Meta-Psychology*. preprint.
1505 <https://doi.org/10.1101/383935>
- 1506 Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2013). Life after p-hacking. Conference
1507 Proceedings. <https://doi.org/10.2139/ssrn.2205186>
- 1508 Stanislaw, H., & Todorov, N. (1999). Calculation of signal detection theory measures.
1509 *Behavior Research Methods, Instruments, & Computers*, 31(1), 137–149. Journal
1510 Article. <https://doi.org/10.3758/BF03207704>
- 1511 Sui, J., He, X., & Humphreys, G. W. (2012). Perceptual effects of social salience: Evidence
1512 from self-prioritization effects on perceptual matching. *Journal of Experimental
1513 Psychology: Human Perception and Performance*, 38(5), 1105–1117. Journal
1514 Article. <https://doi.org/10.1037/a0029792>
- 1515 Wiecki, T. V., Sofer, I., & Frank, M. J. (2013). HDDM: Hierarchical bayesian estimation of
1516 the drift-diffusion model in python. *Frontiers in Neuroinformatics*, 7.
1517 <https://doi.org/10.3389/fninf.2013.00014>

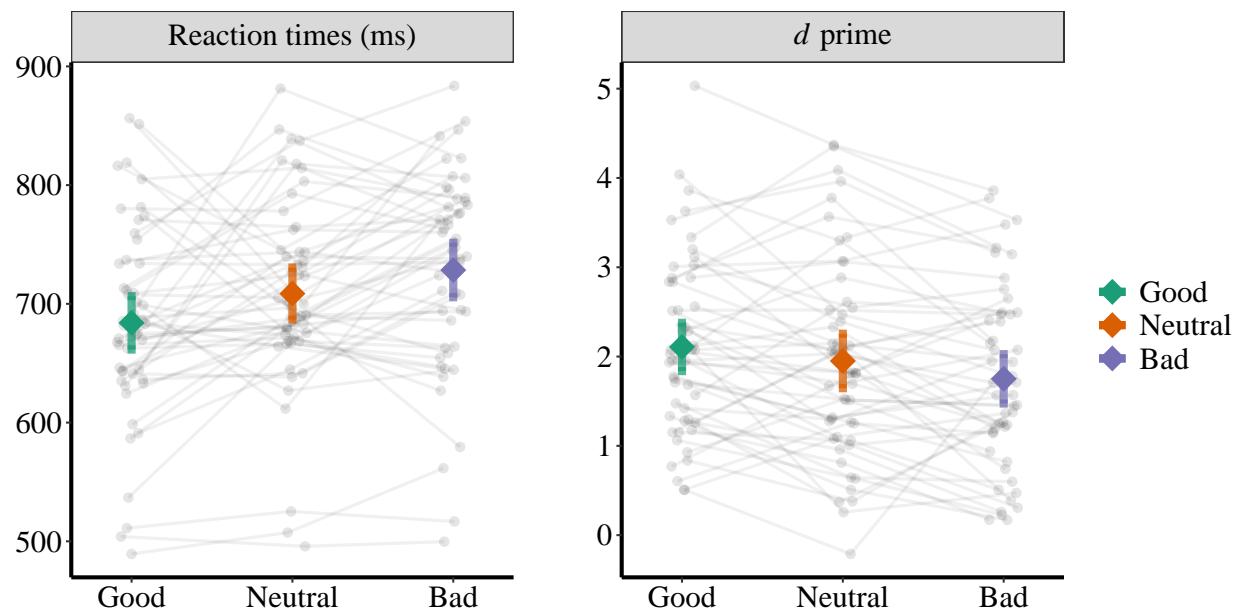


Figure 1. RT and d prime of Experiment 1a.

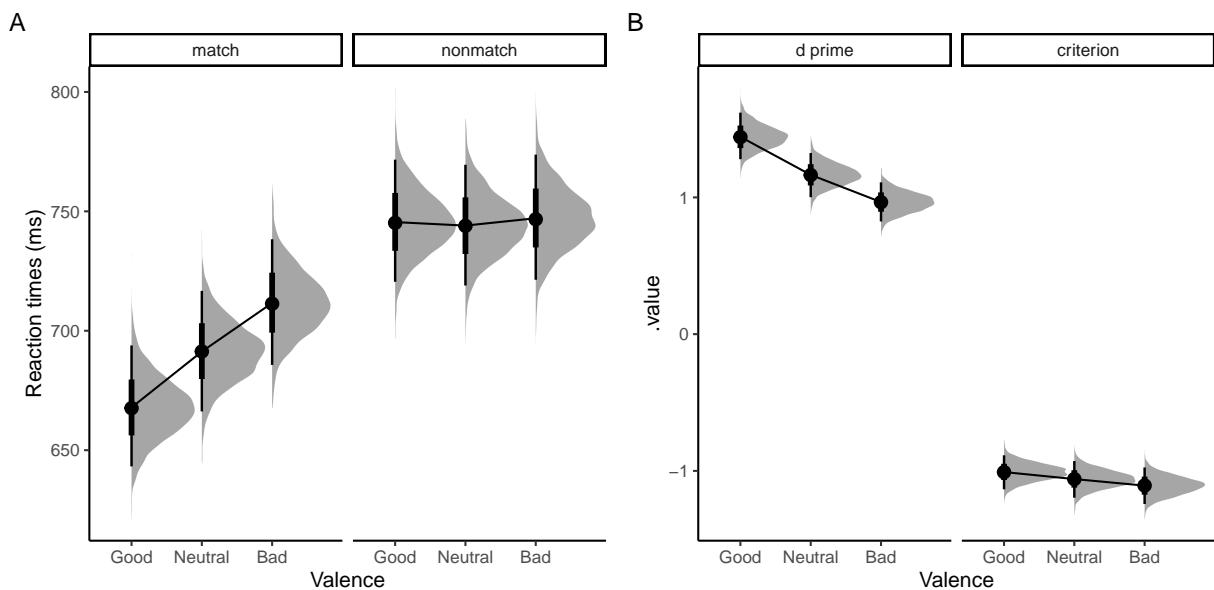


Figure 2. Exp1a: Results of Bayesian GLM analysis.

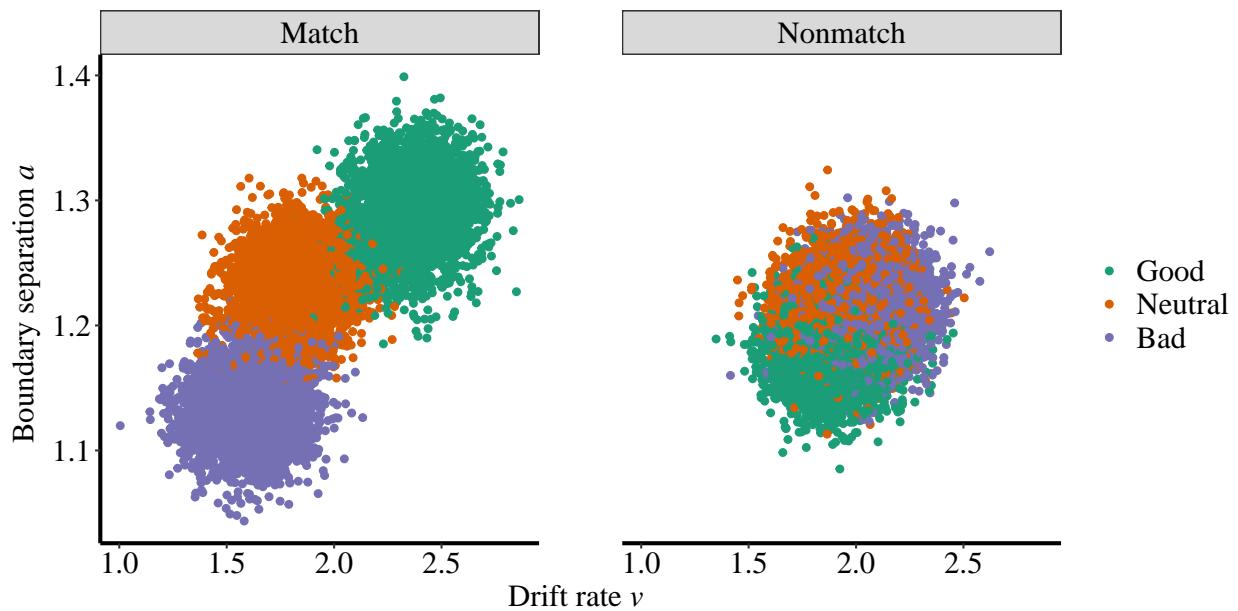


Figure 3. Exp1a: Results of HDDM.

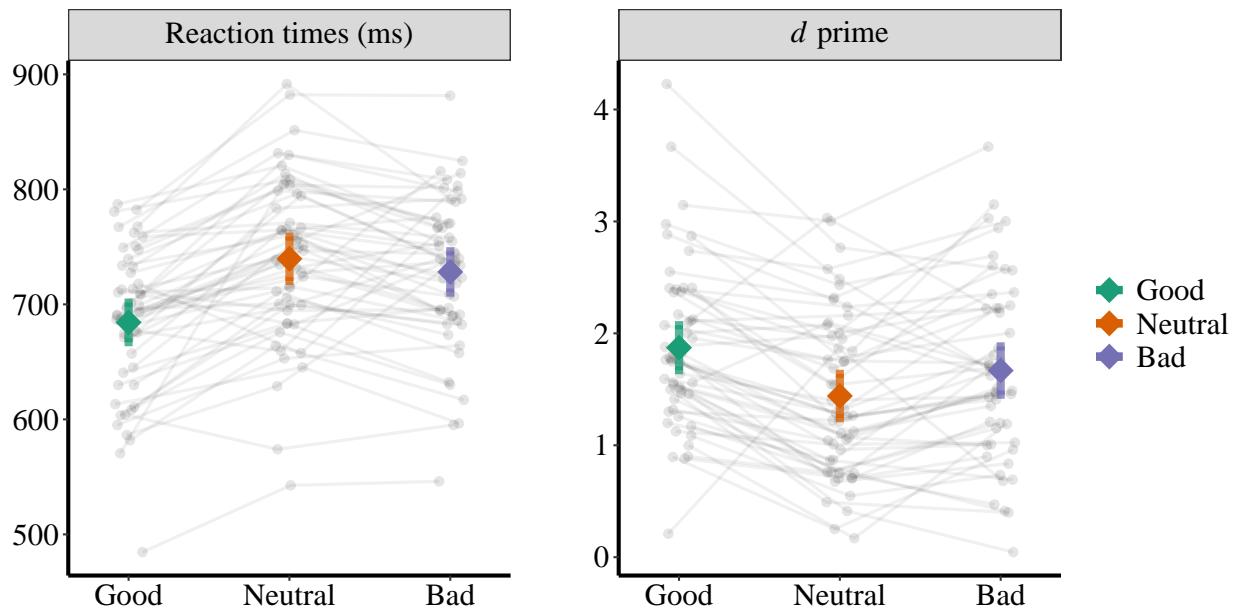


Figure 4. RT and d' of Experiment 1b.

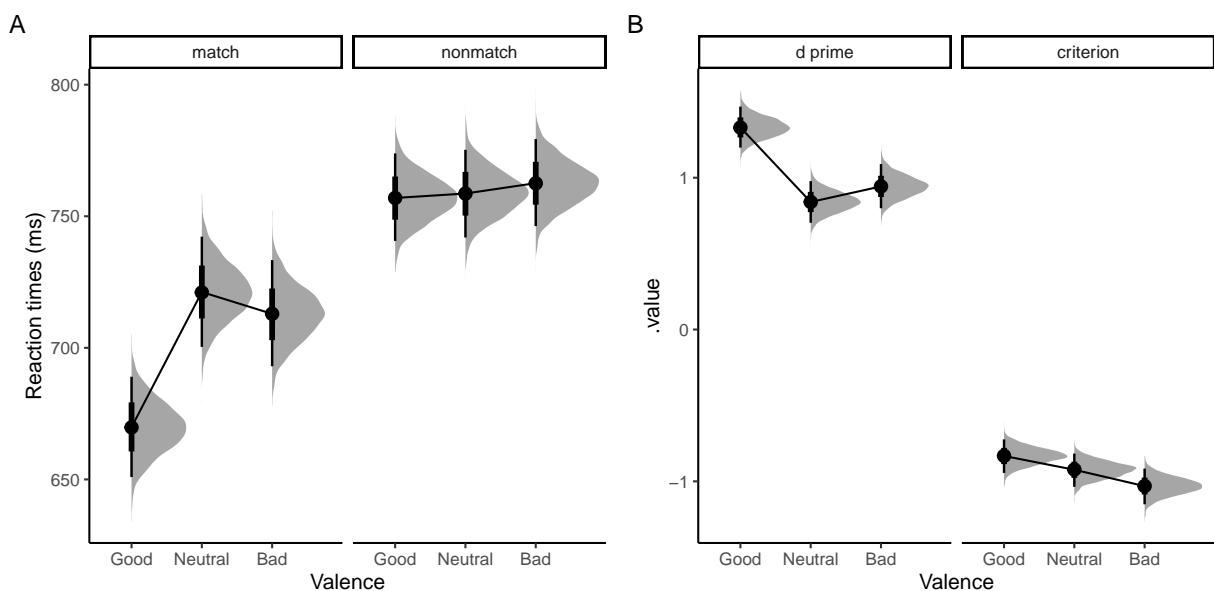


Figure 5. Exp1b: Results of Bayesian GLM analysis.

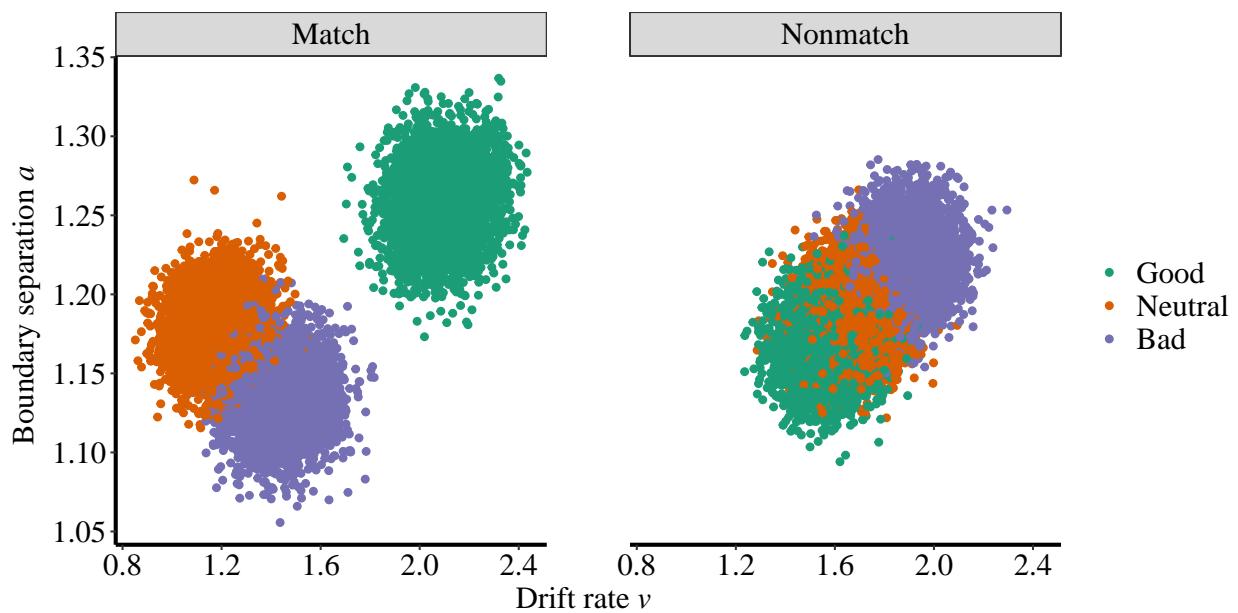


Figure 6. Exp1b: Results of HDDM.

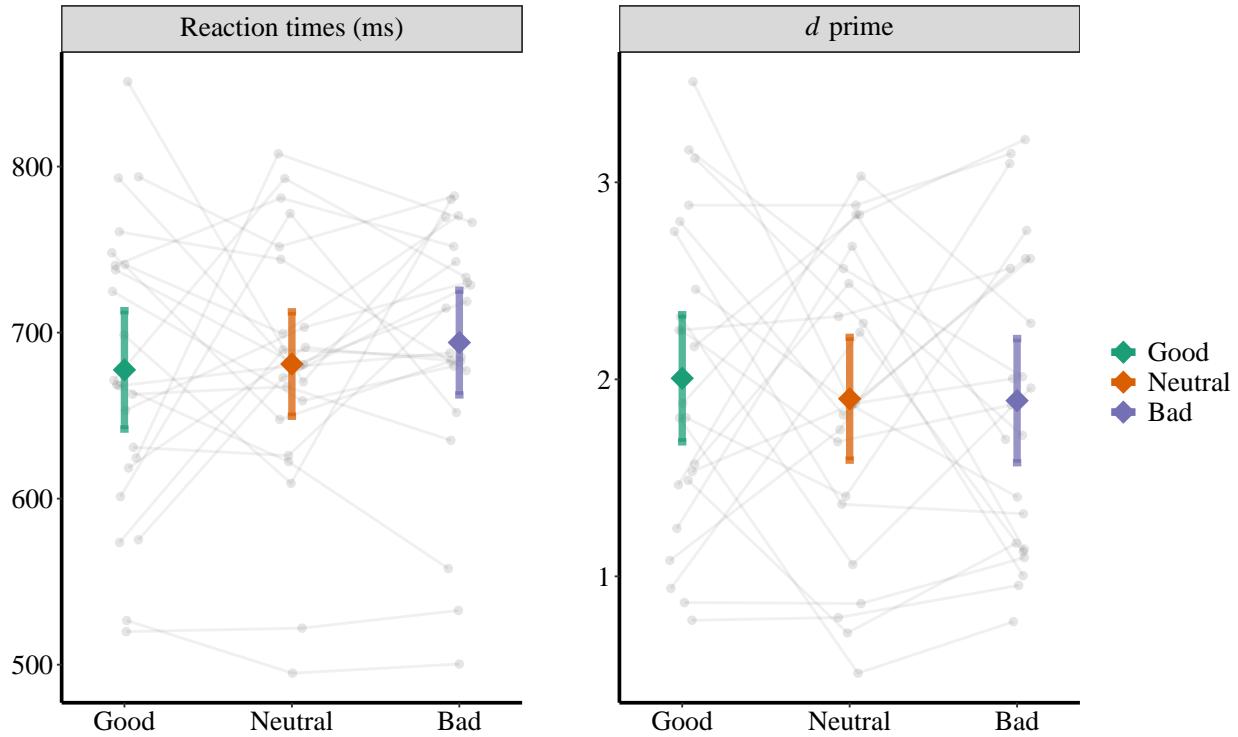


Figure 7. RT and d' prime of Experiment 1c.

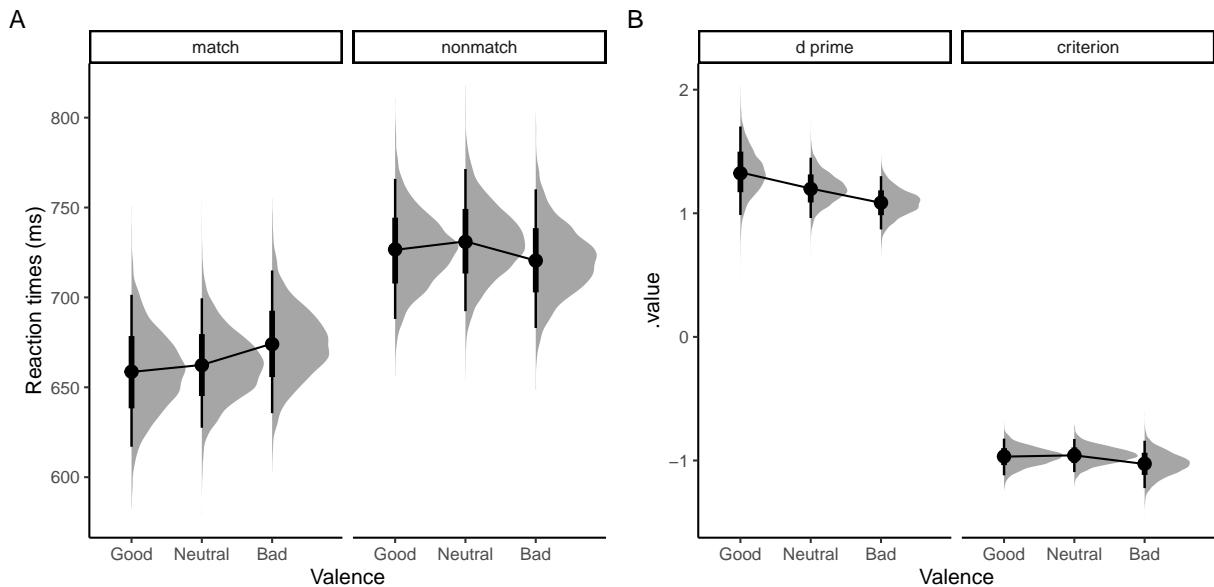


Figure 8. Exp1c: Results of Bayesian GLM analysis.

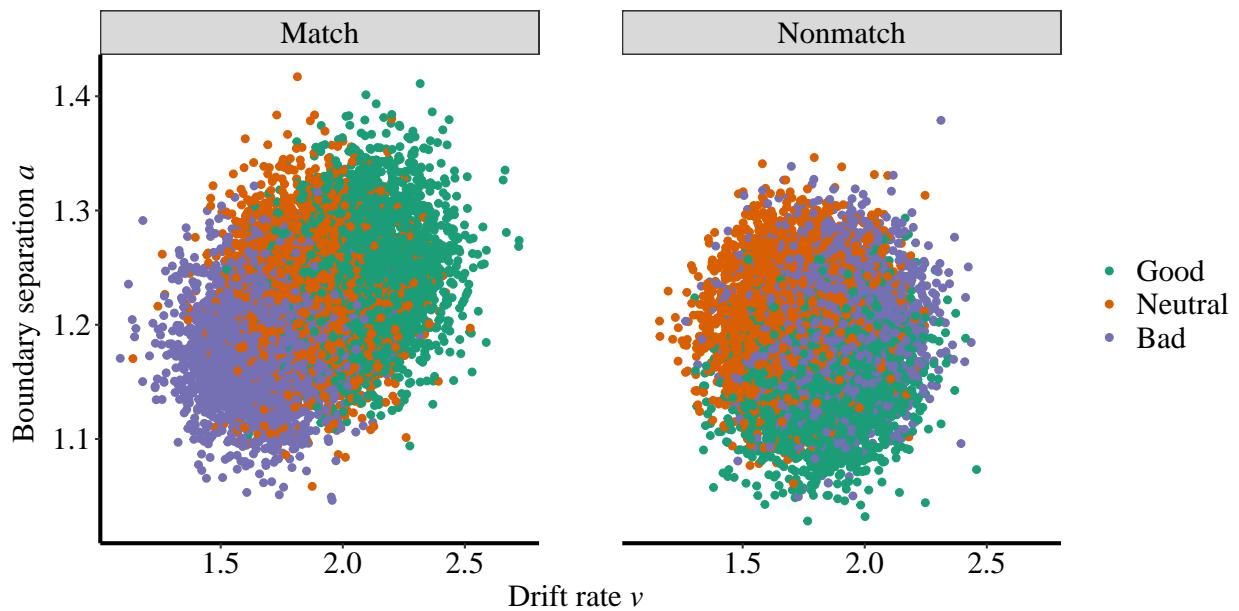


Figure 9. Exp1c: Results of HDDM.

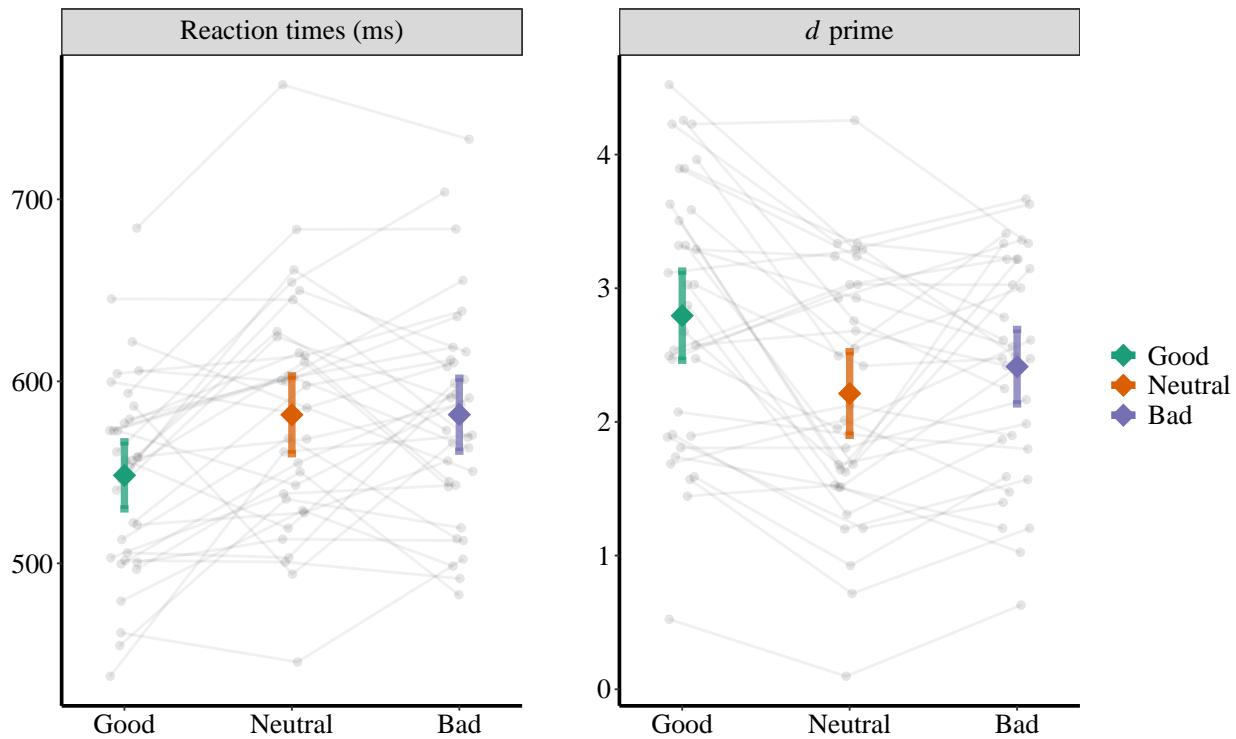


Figure 10. RT and d' of Experiment 2.

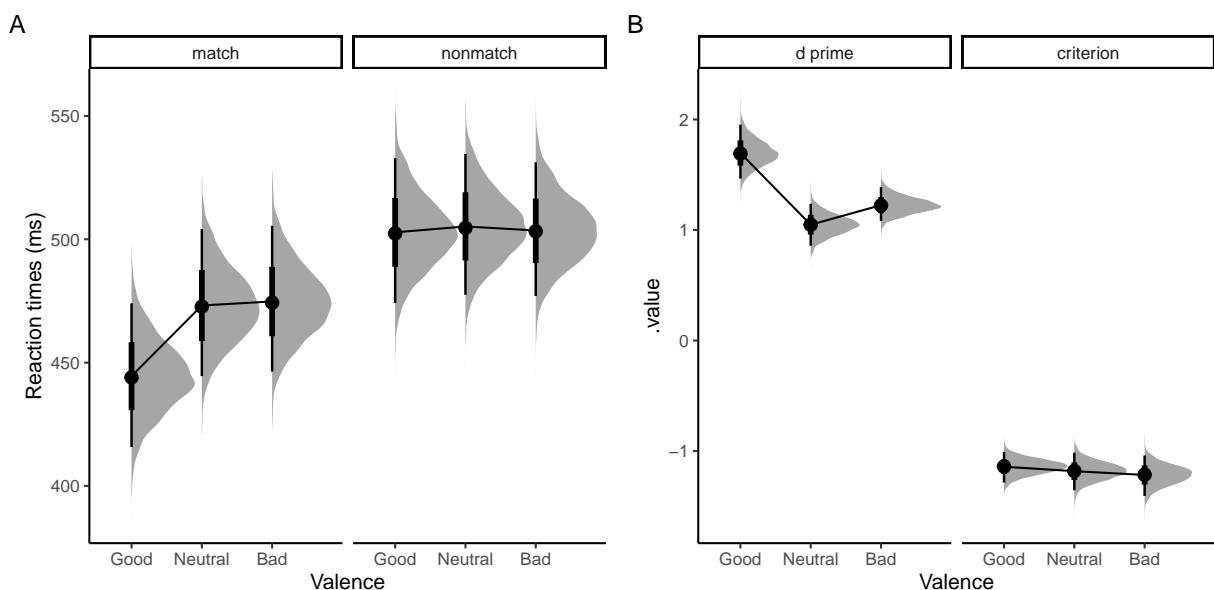


Figure 11. Exp2: Results of Bayesian GLM analysis.

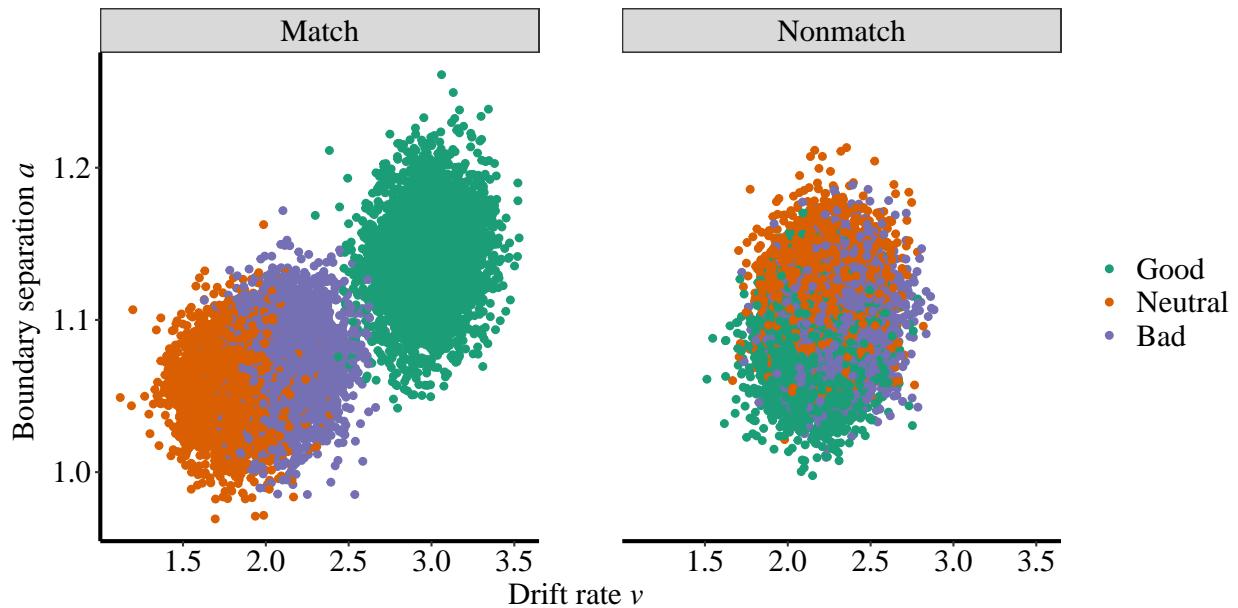


Figure 12. Exp2: Results of HDDM.

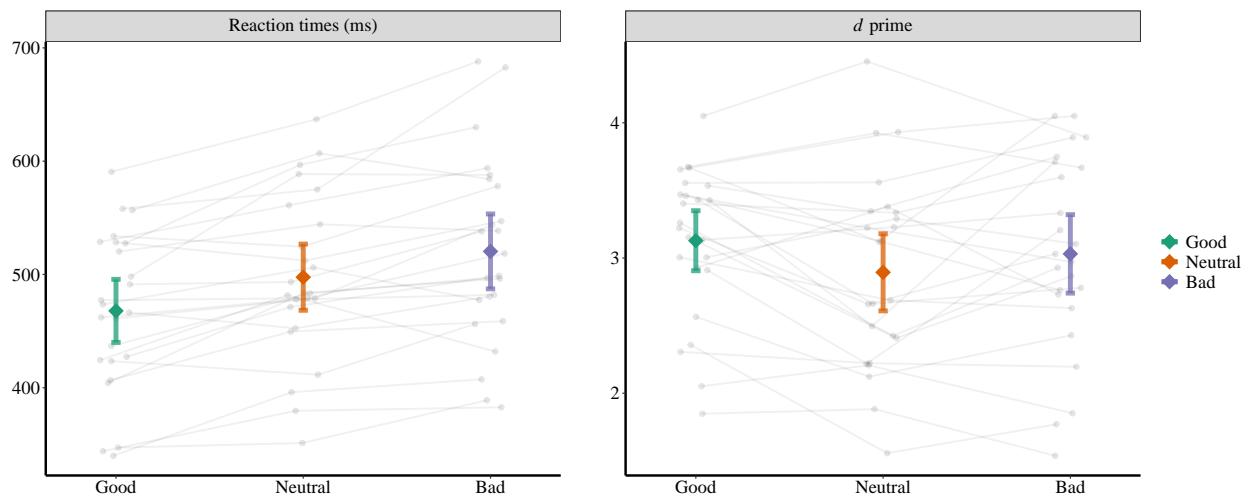


Figure 13. RT and d' prime of Experiment 6a.

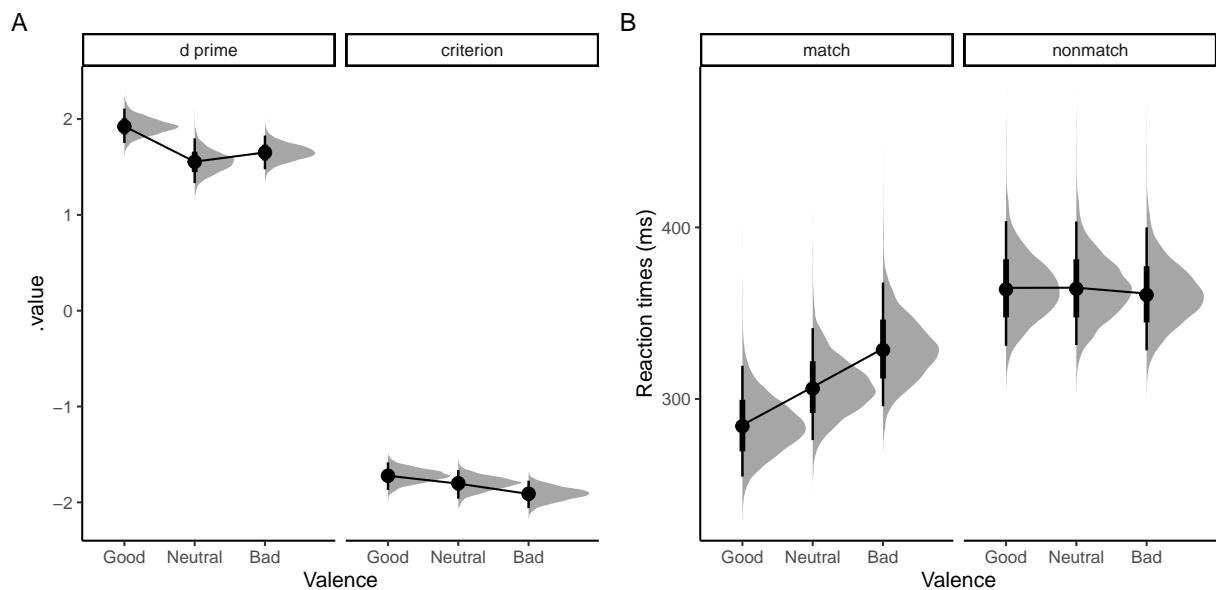


Figure 14. Exp6a: Results of Bayesian GLM analysis.

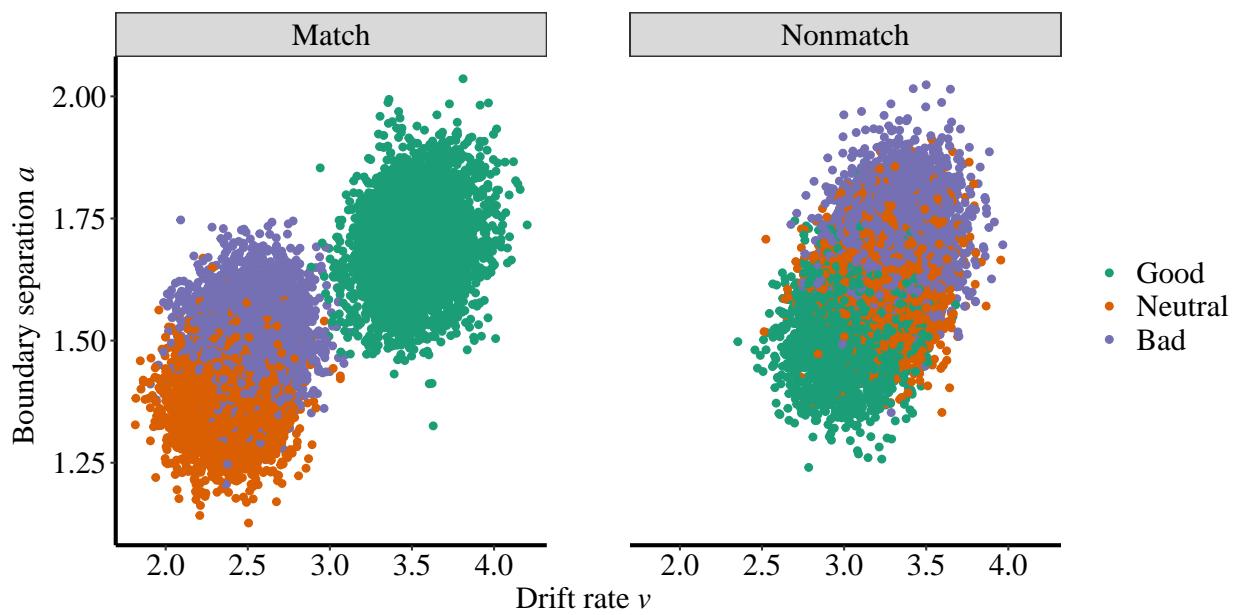


Figure 15. exp6a: Results of HDDM.

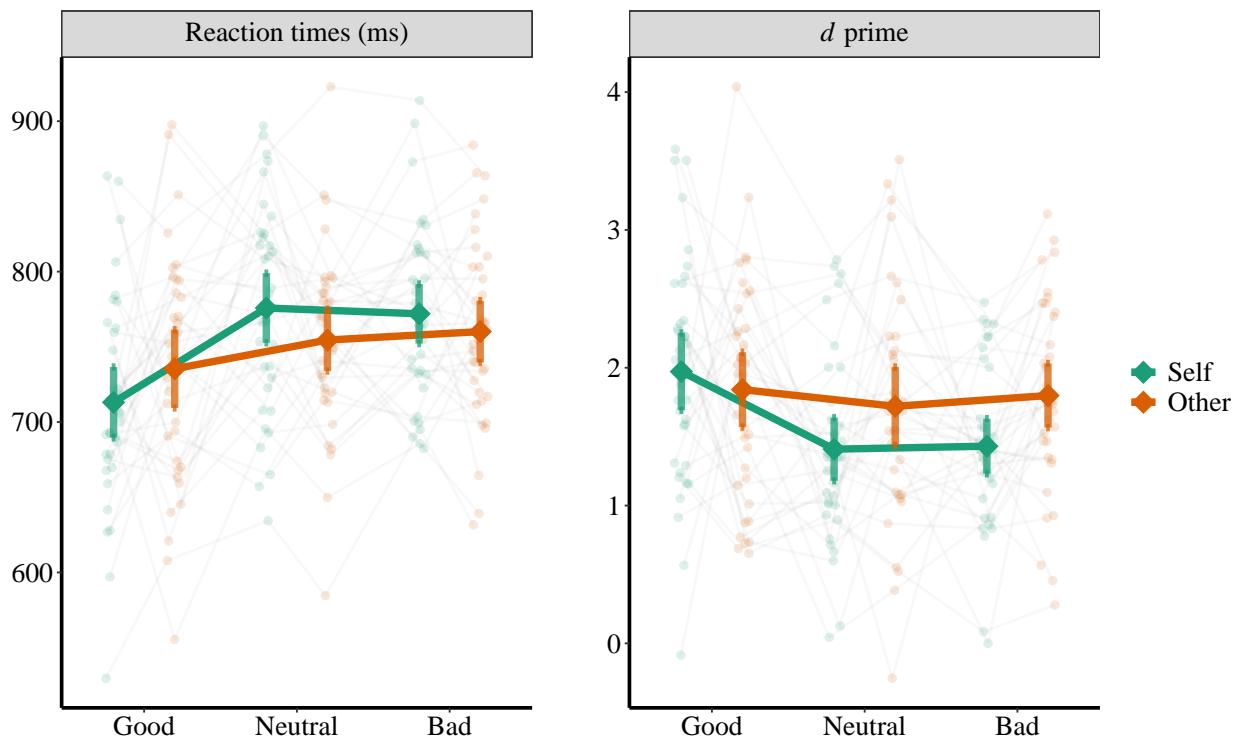


Figure 16. RT and d prime of Experiment 3a.

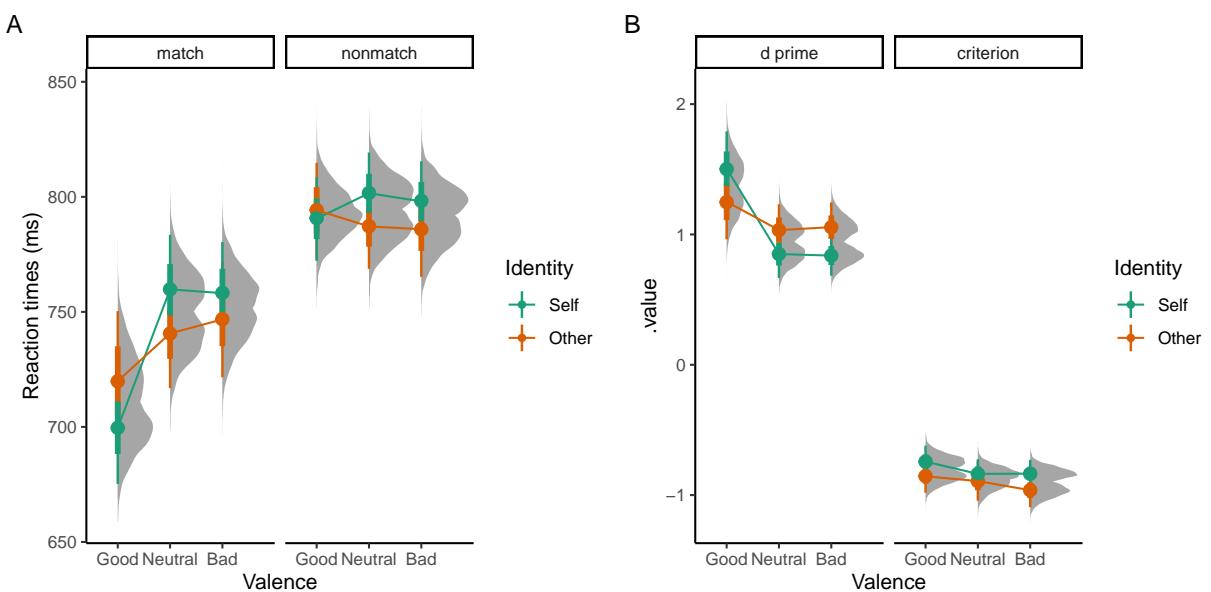


Figure 17. Exp3a: Results of Bayesian GLM analysis.

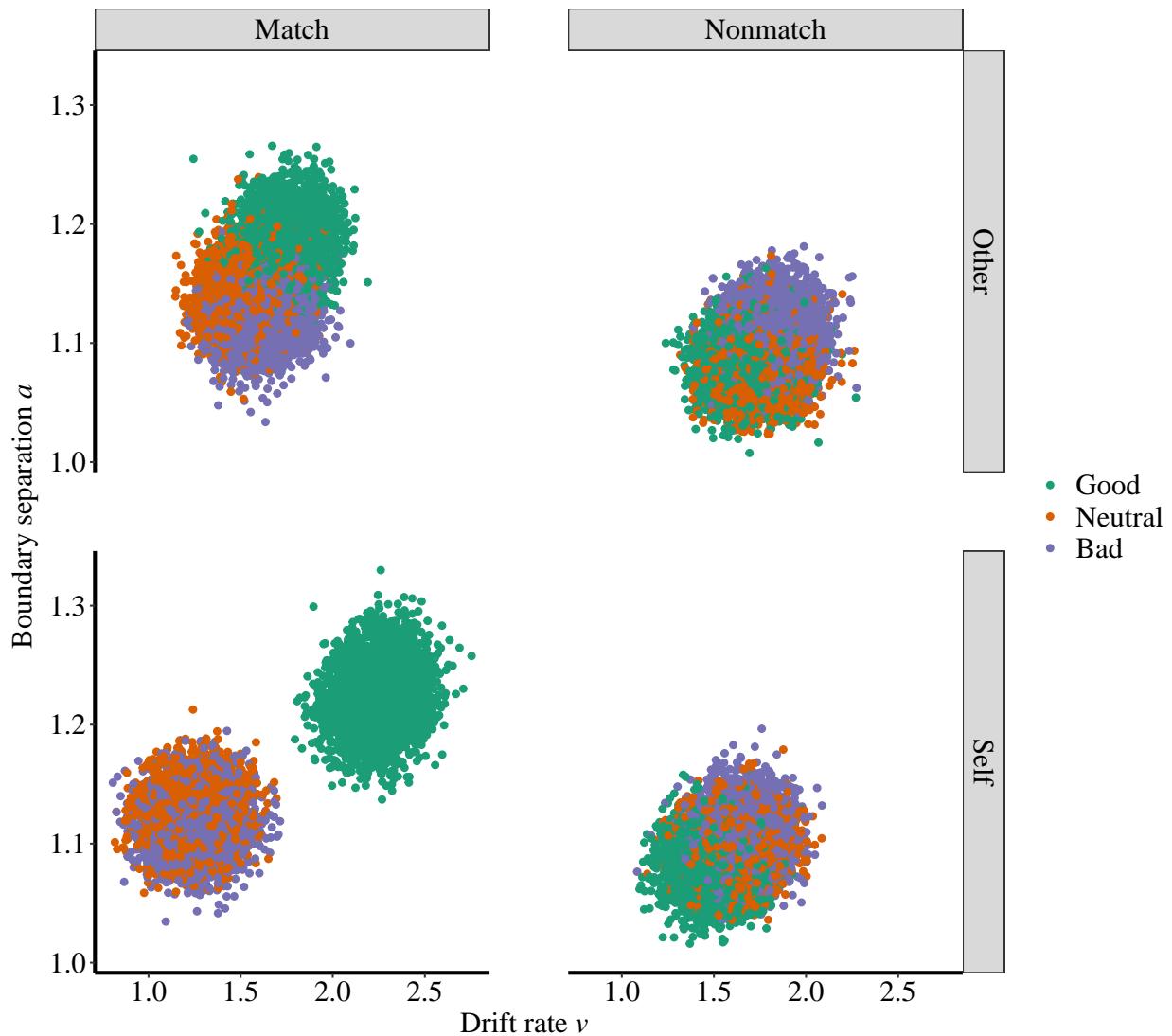


Figure 18. Exp3a: Results of HDDM.

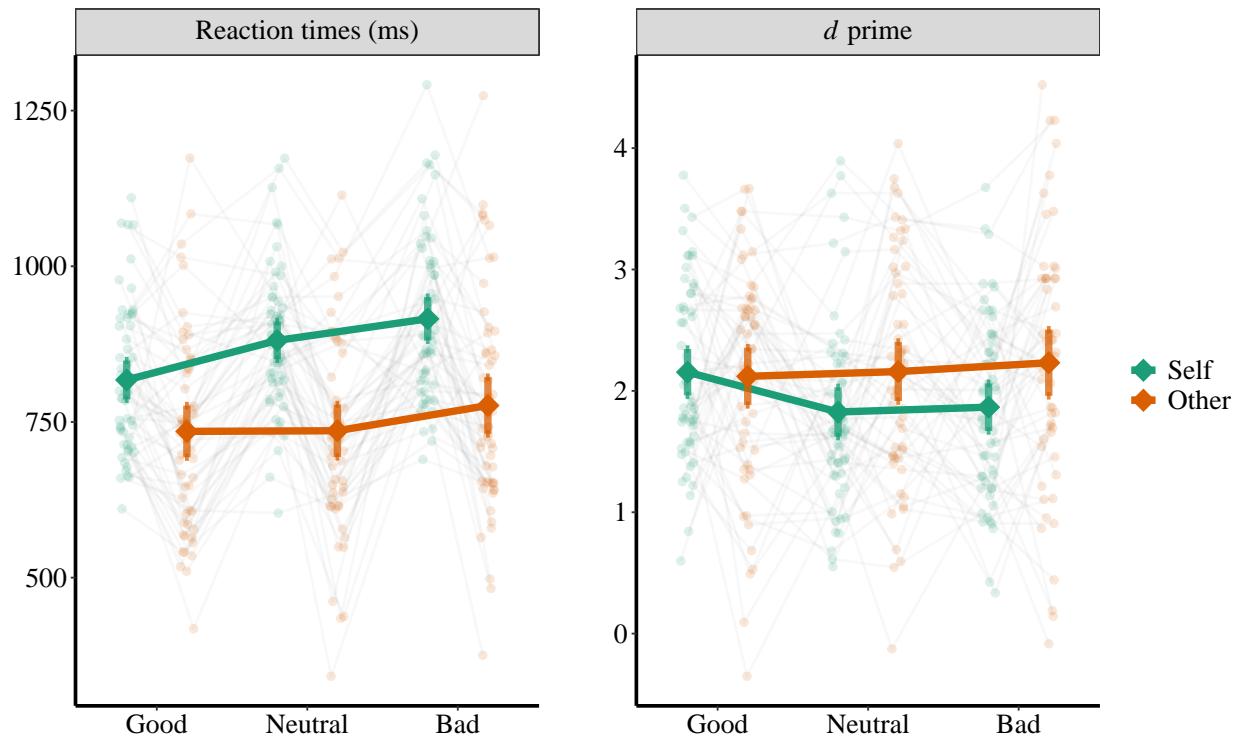


Figure 19. RT and d' prime of Experiment 3b.

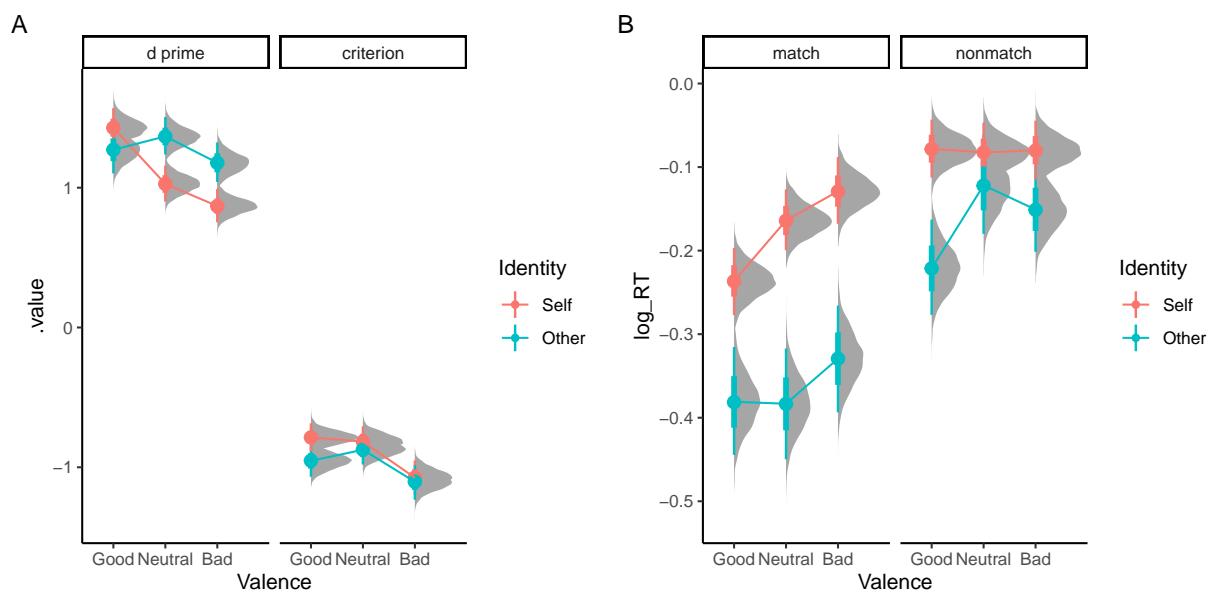


Figure 20. exp3b: Results of Bayesian GLM analysis.

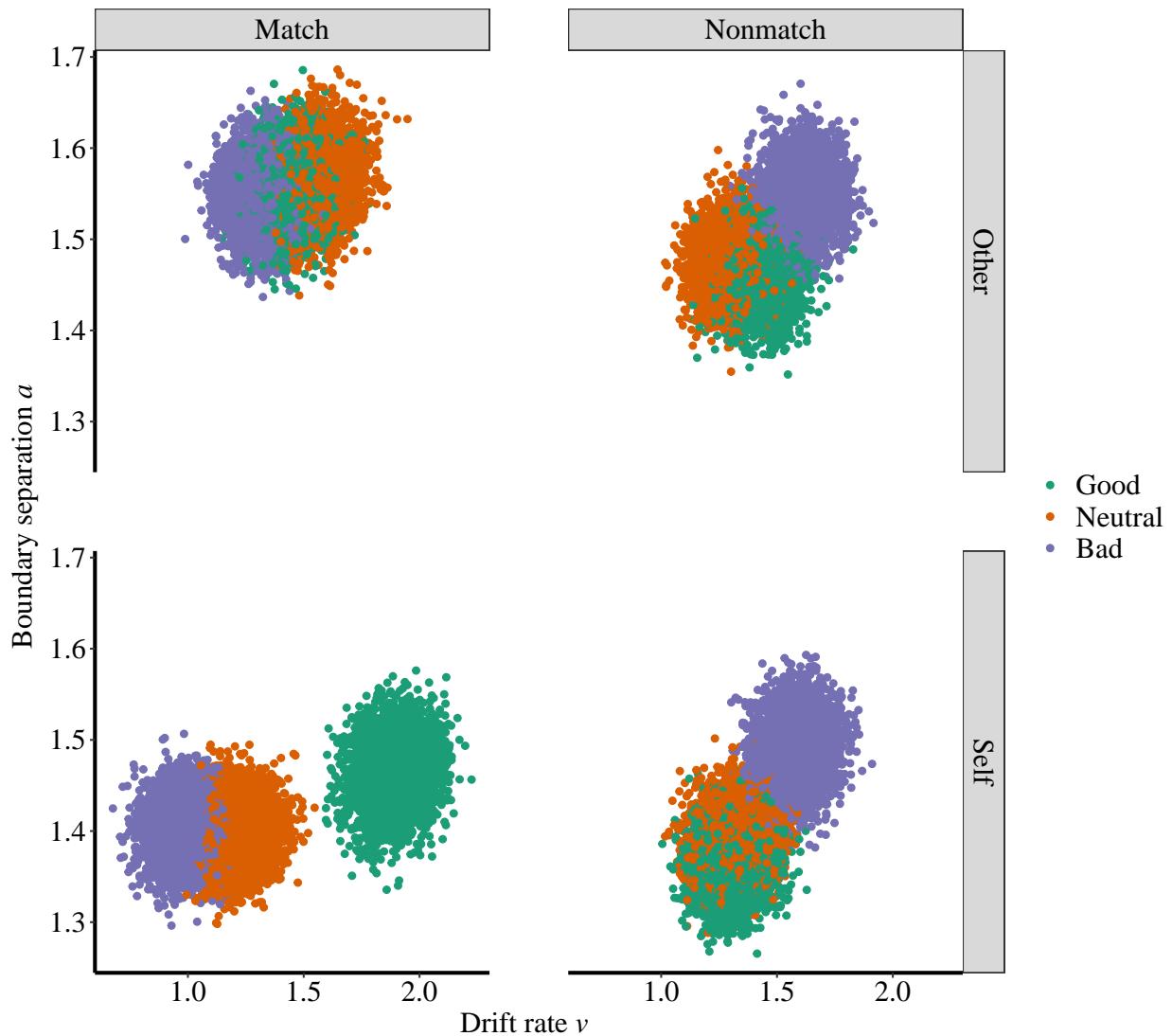


Figure 21. exp3b: Results of HDDM.

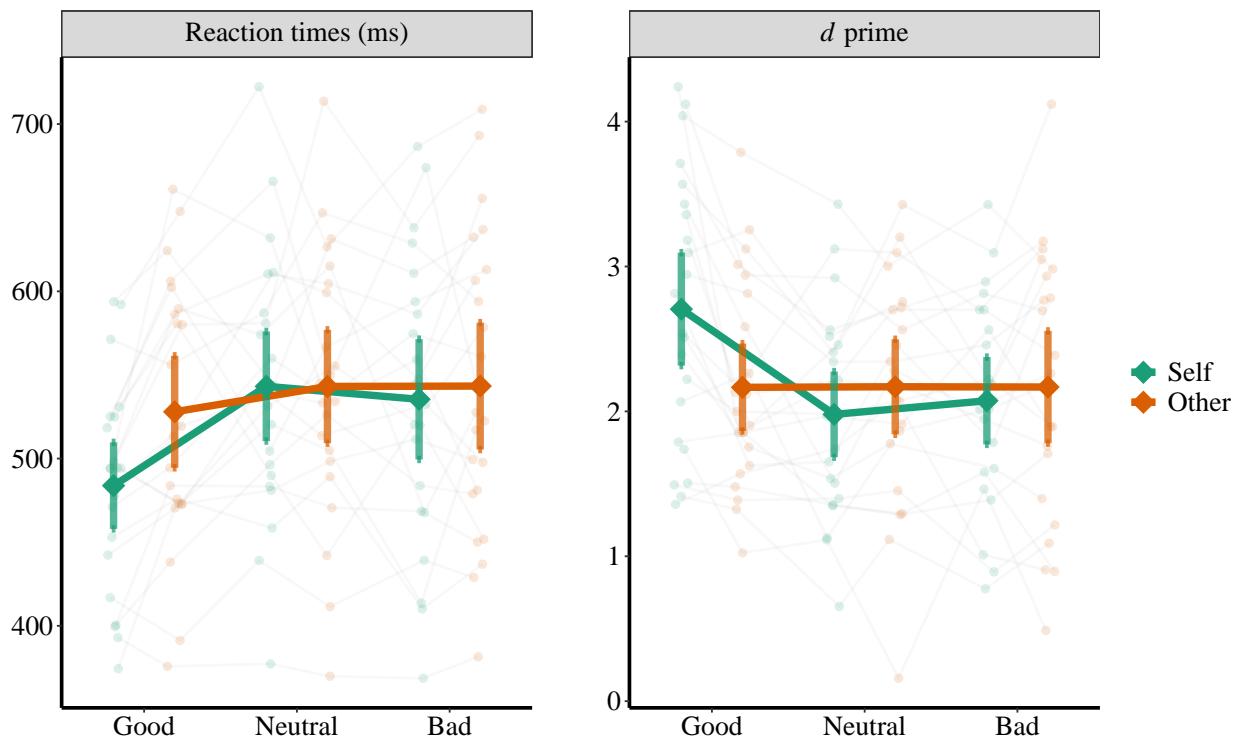


Figure 22. RT and d prime of Experiment 6b.

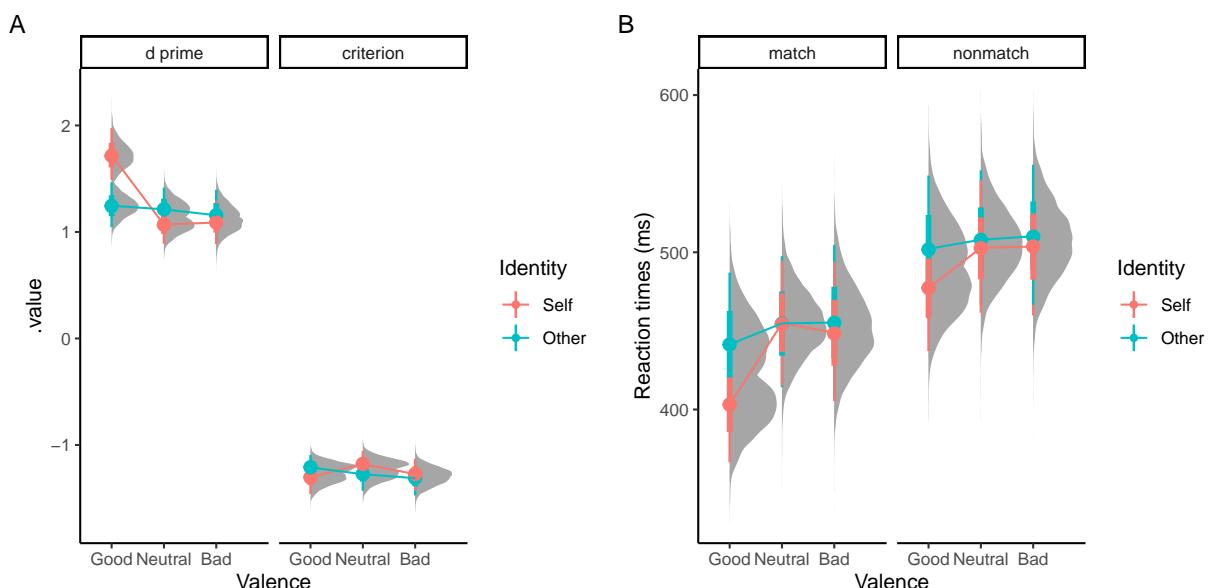


Figure 23. exp6b_d1: Results of Bayesian GLM analysis.

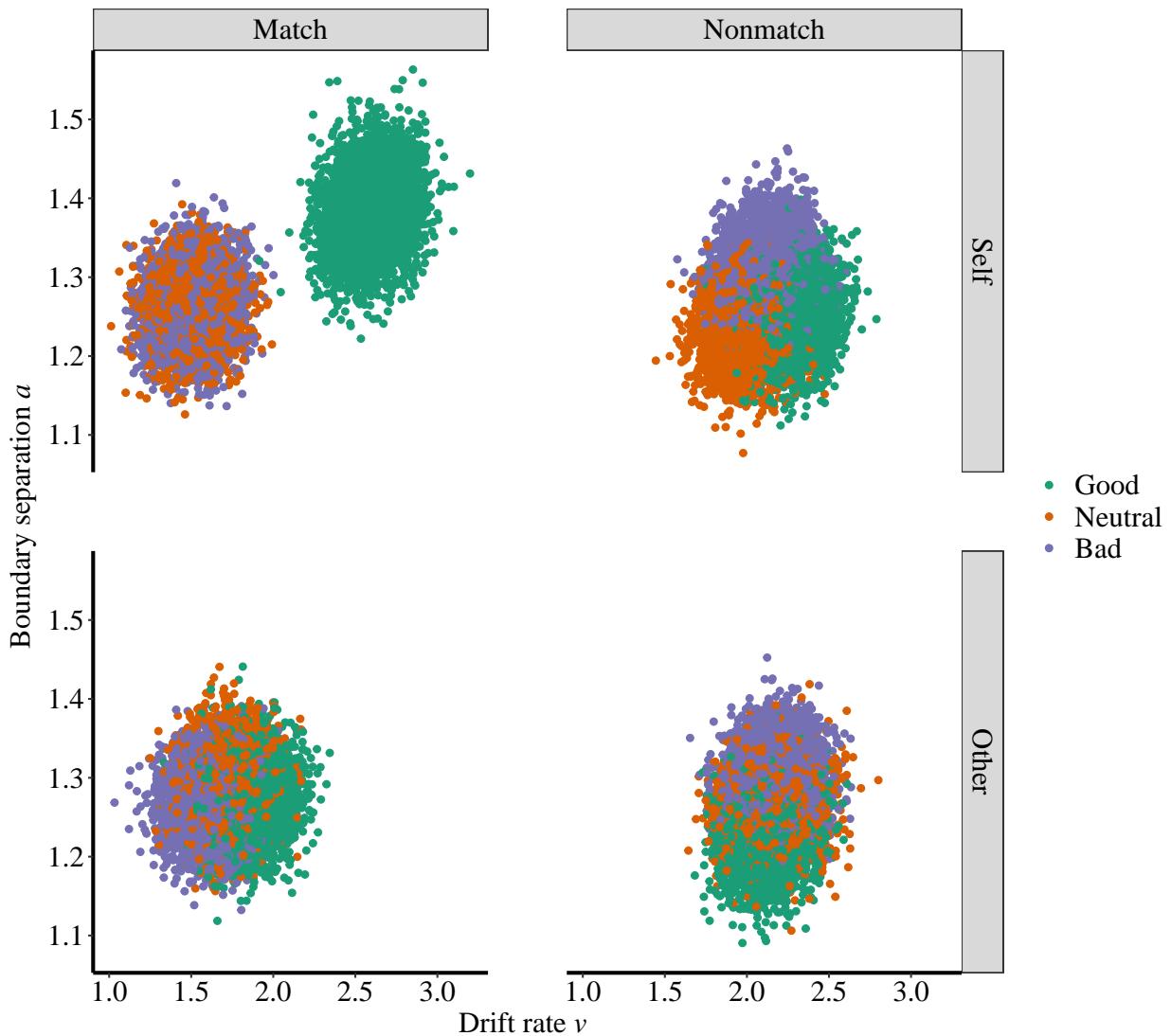


Figure 24. exp6b: Results of HDDM (Day 1).

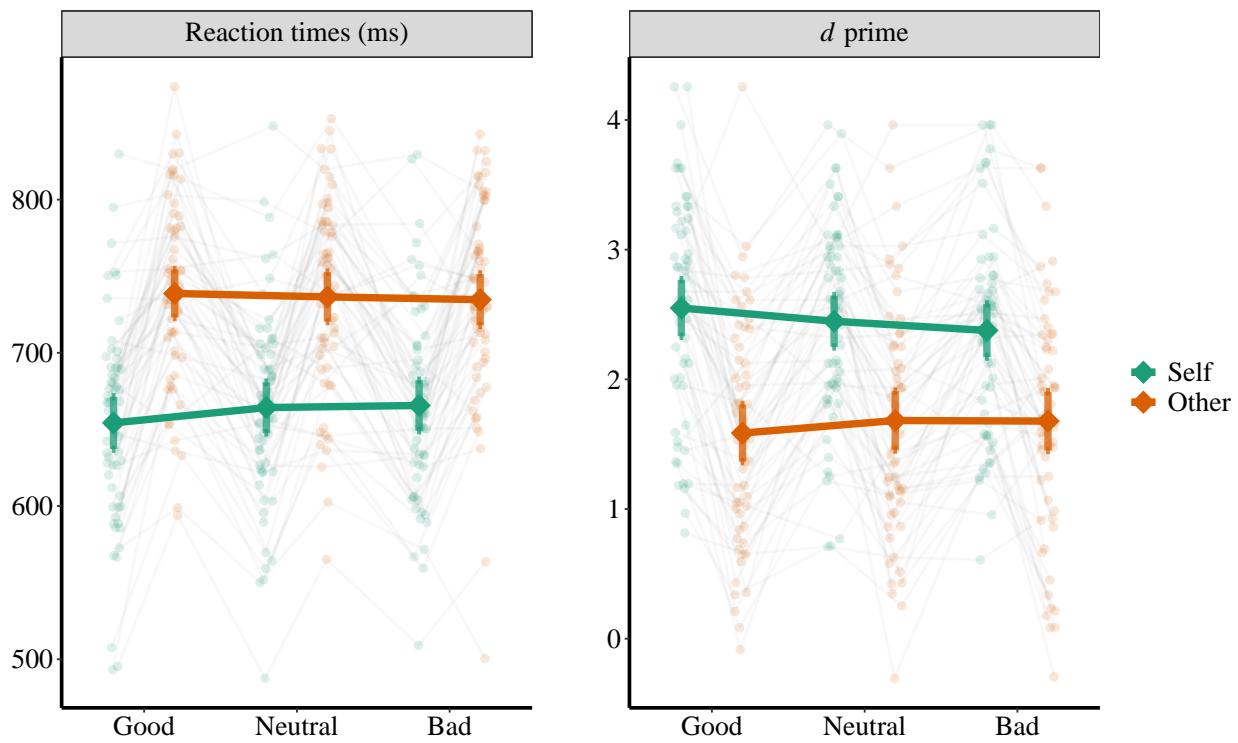
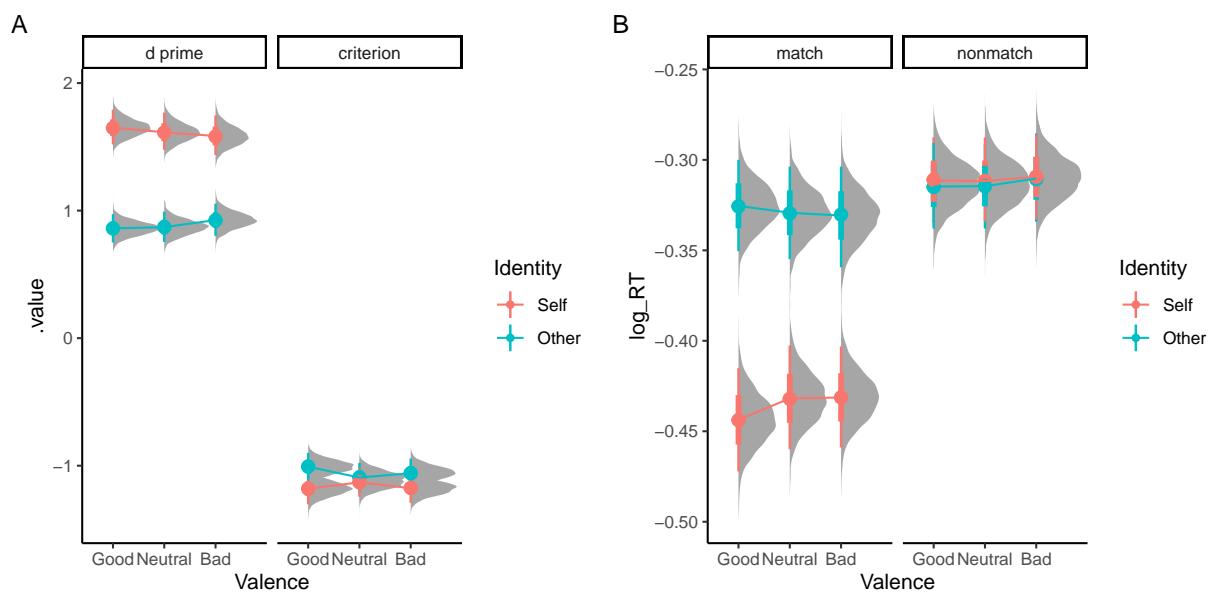
Figure 25. RT and d' of Experiment 4a.

Figure 26. exp4a: Results of Bayesian GLM analysis.

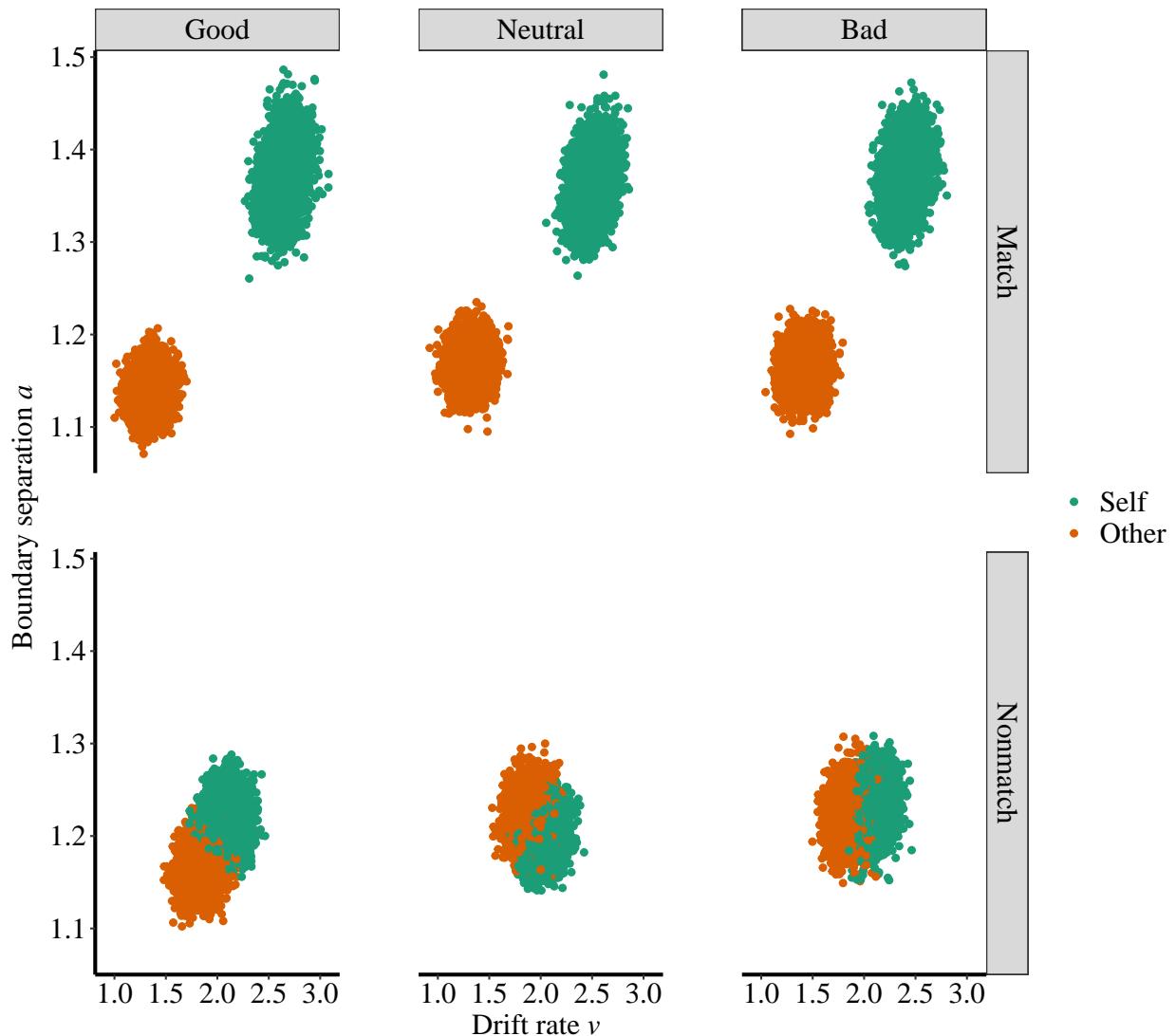


Figure 27. exp4a: Results of HDDM.

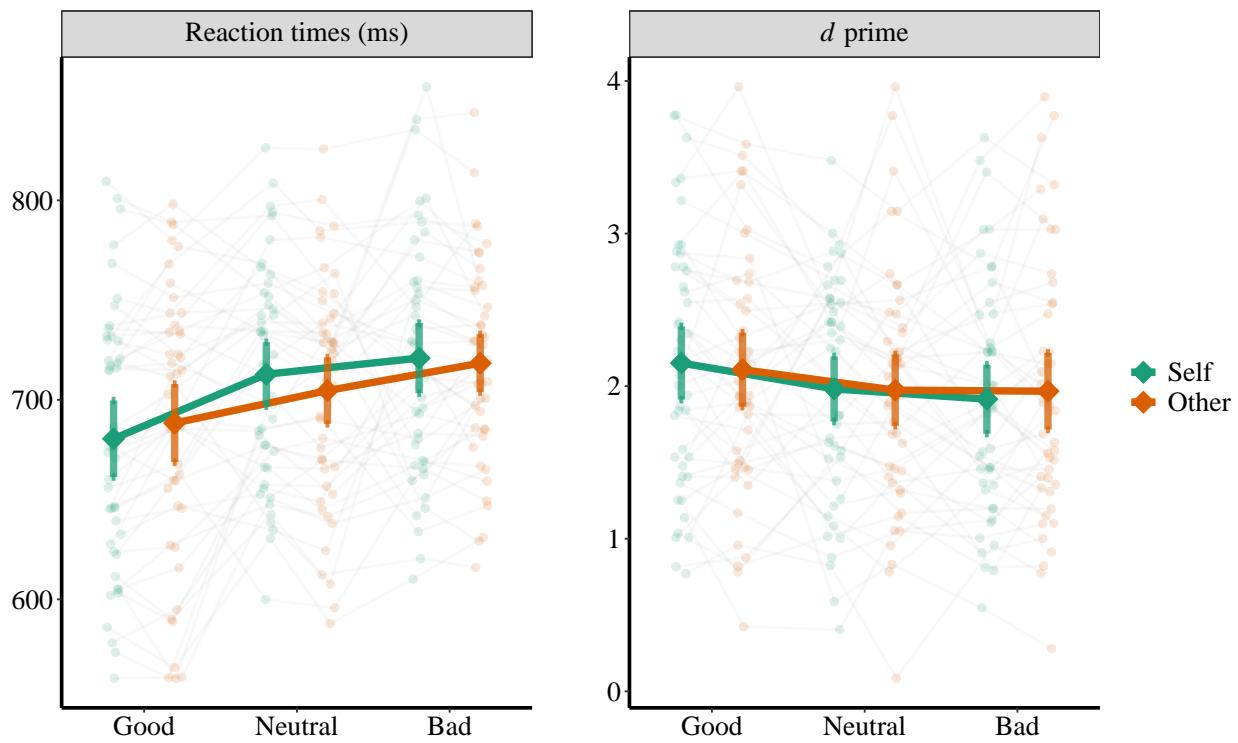


Figure 28. RT and d' of Experiment 4b.

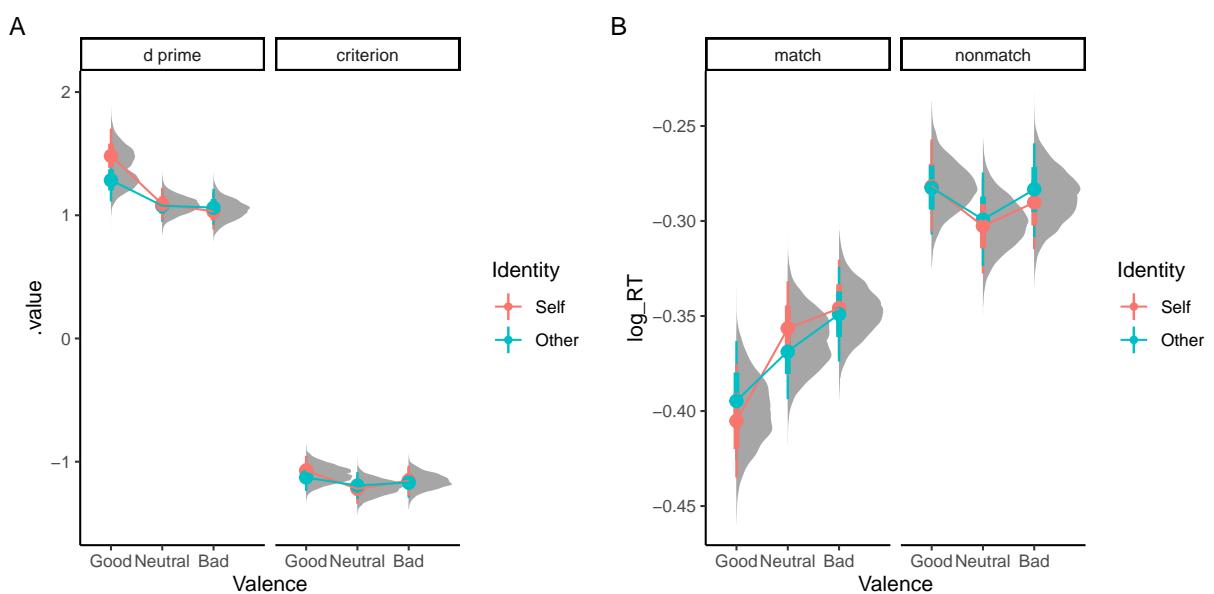


Figure 29. exp4b: Results of Bayesian GLM analysis.

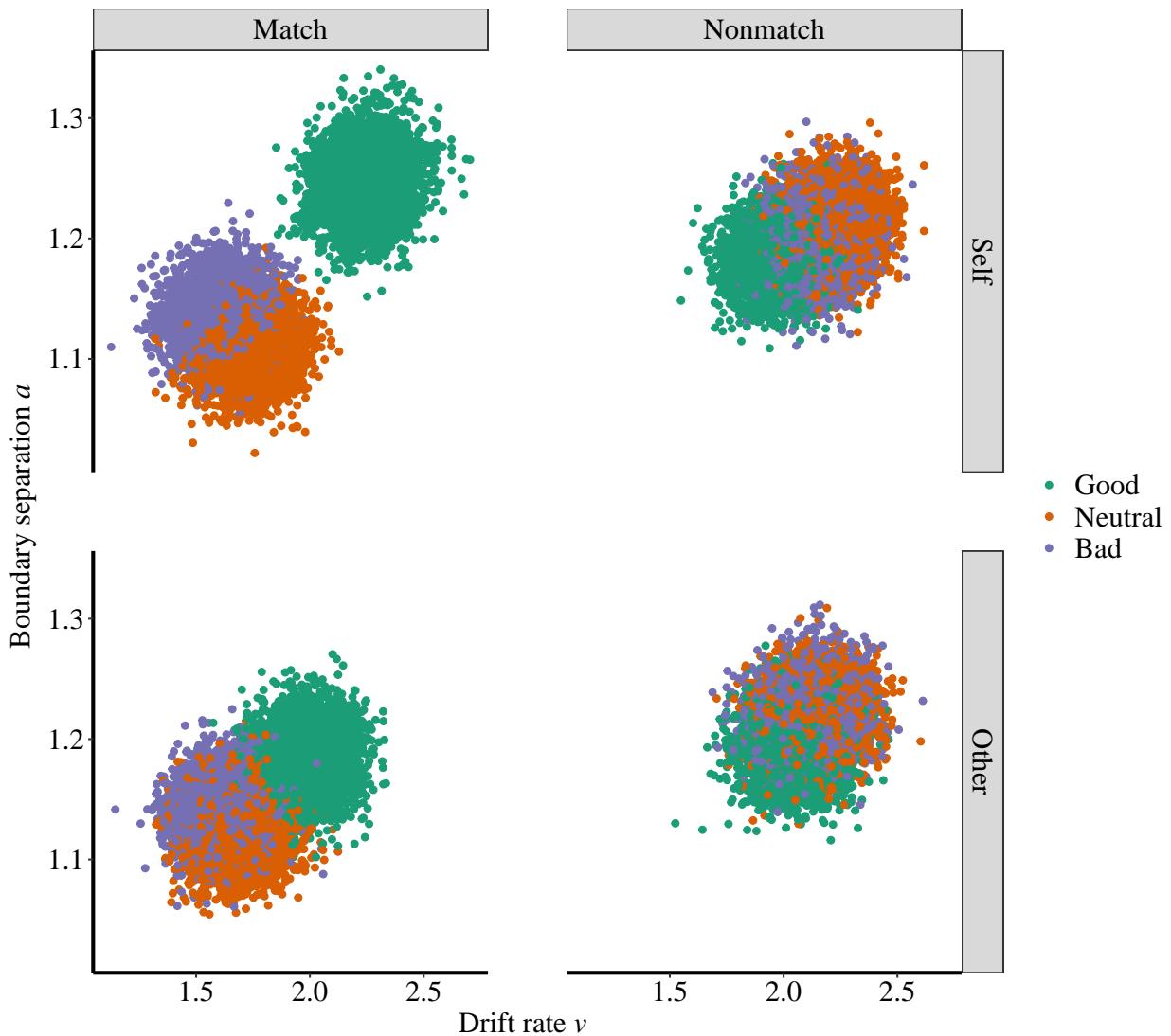


Figure 30. exp4b: Results of HDDM.

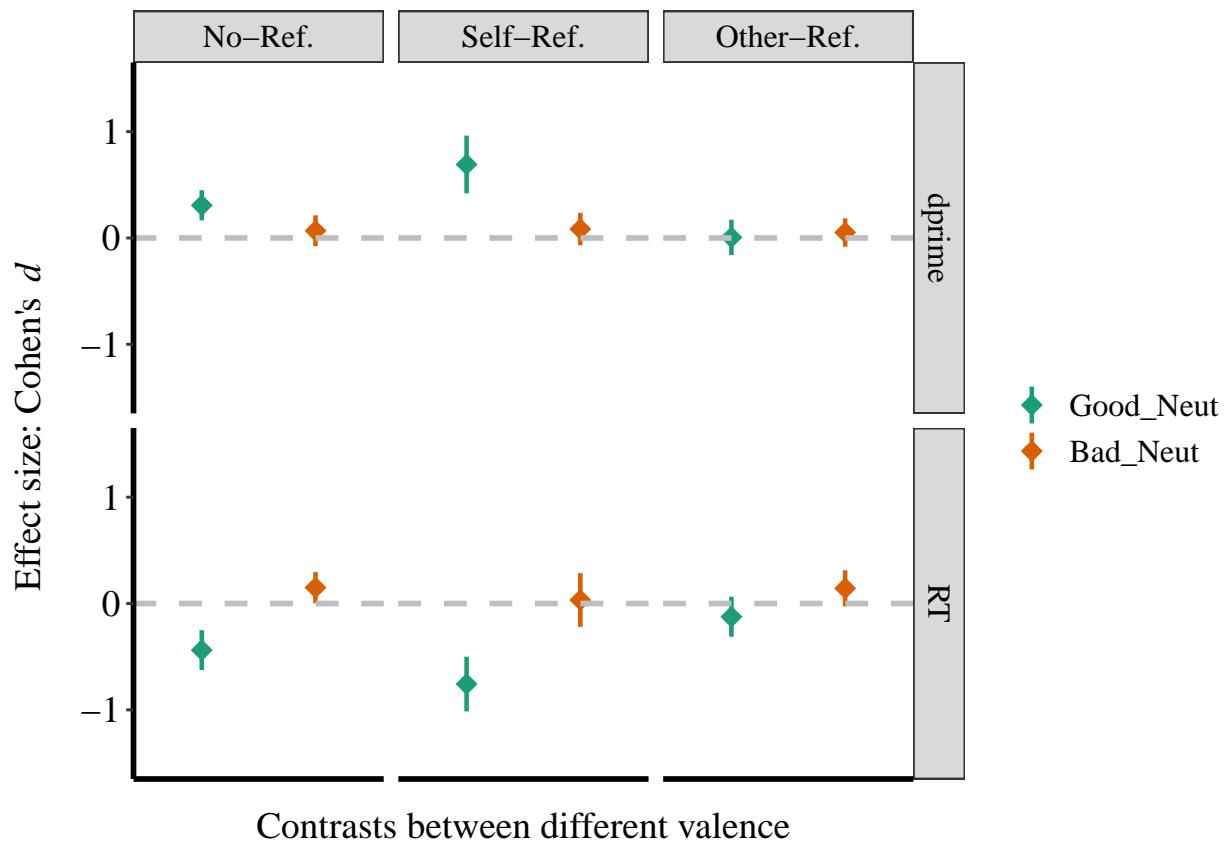


Figure 31. Effect size (Cohen's d) of Valence.

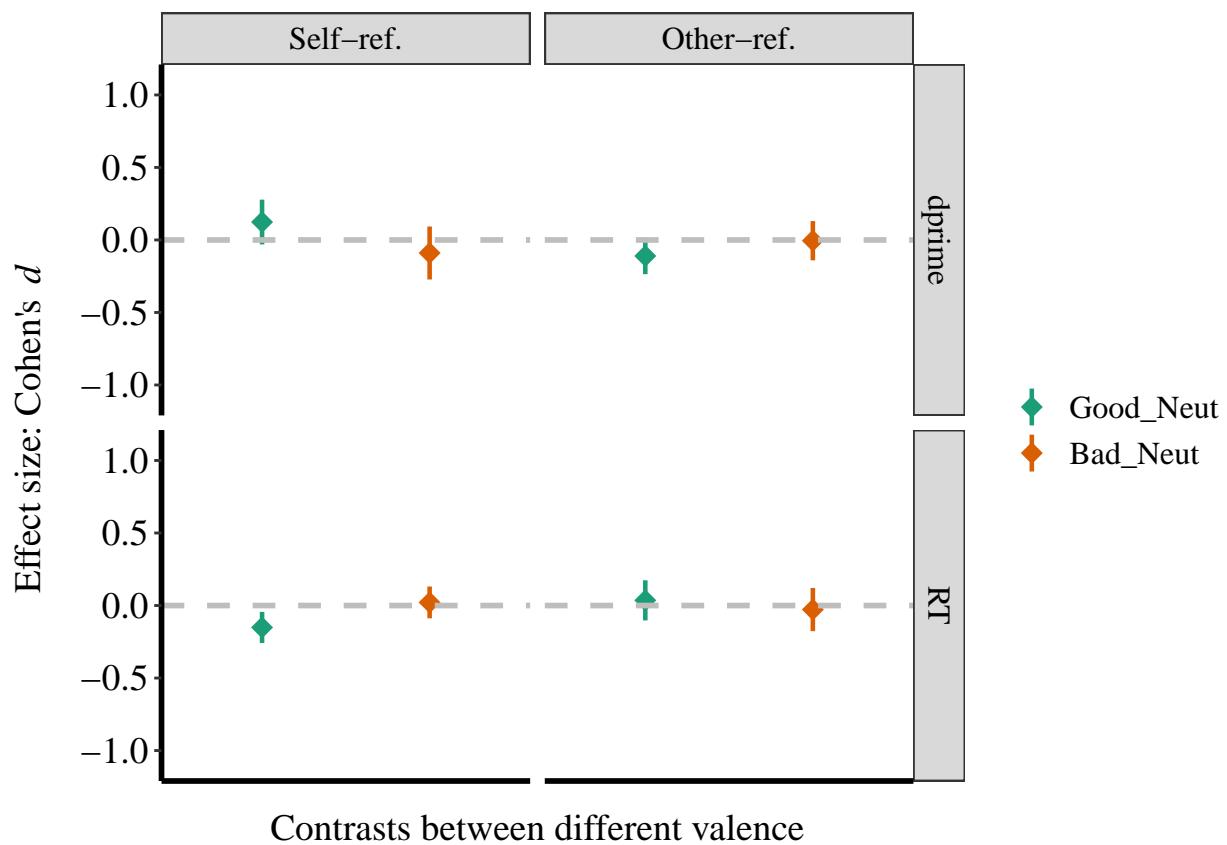


Figure 32. Effect size (Cohen's d) of Valence in Exp4a.

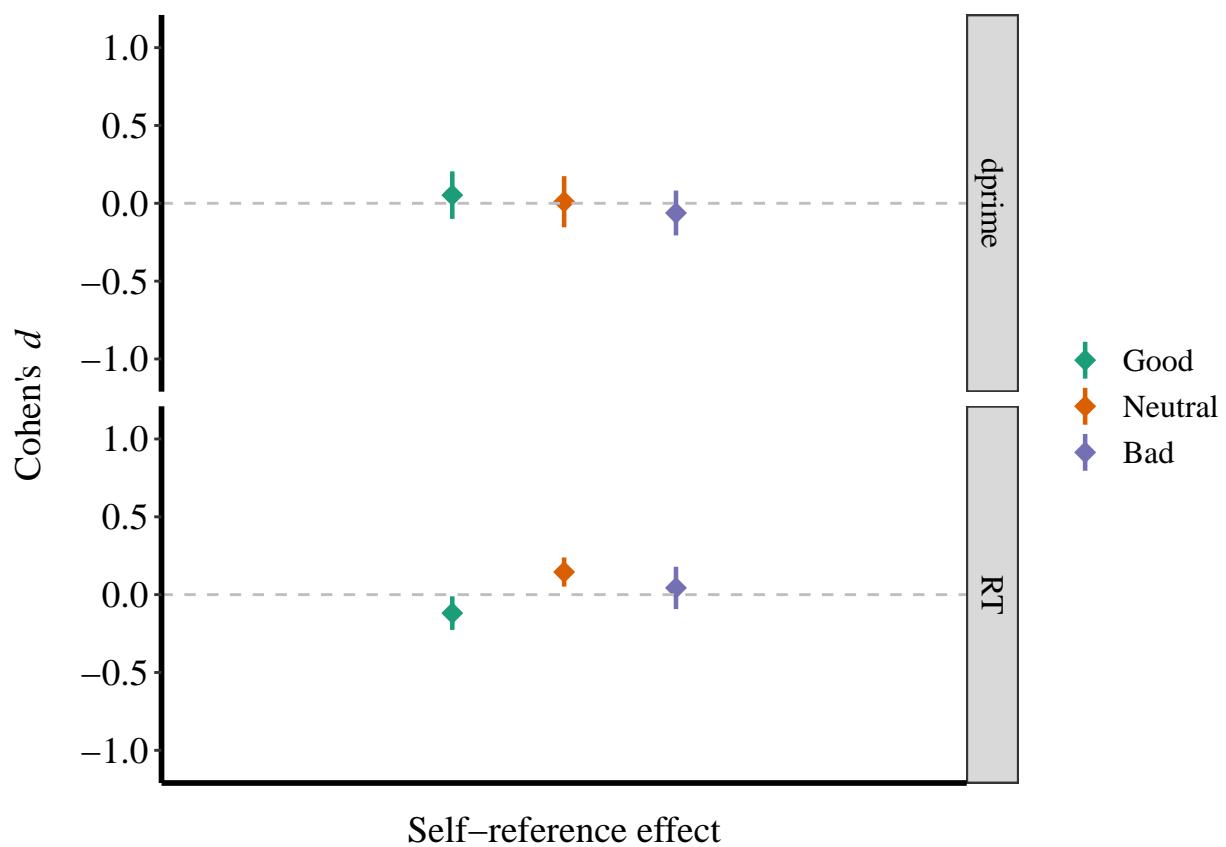


Figure 33. Effect size (Cohen's d) of Valence in Exp4b.

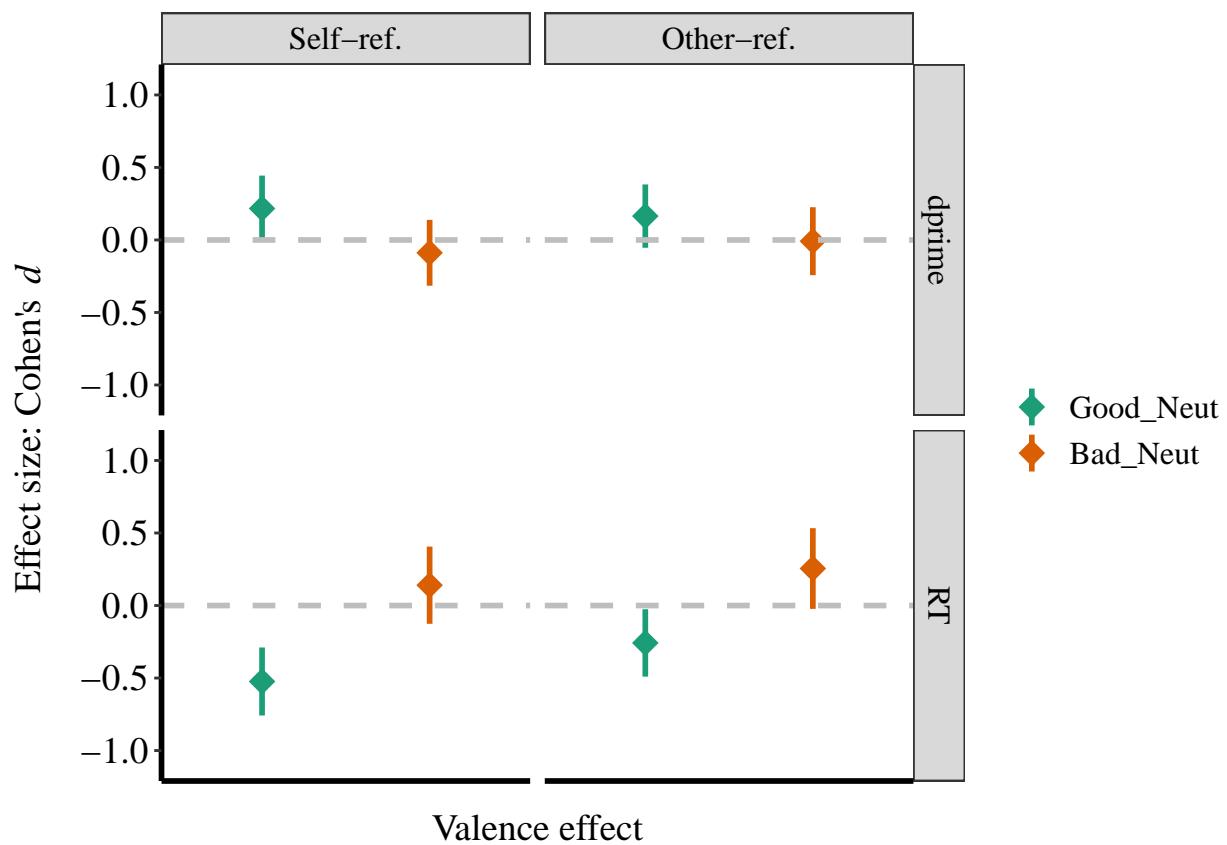


Figure 34. Effect size (Cohen's d) of Valence in Exp4b.

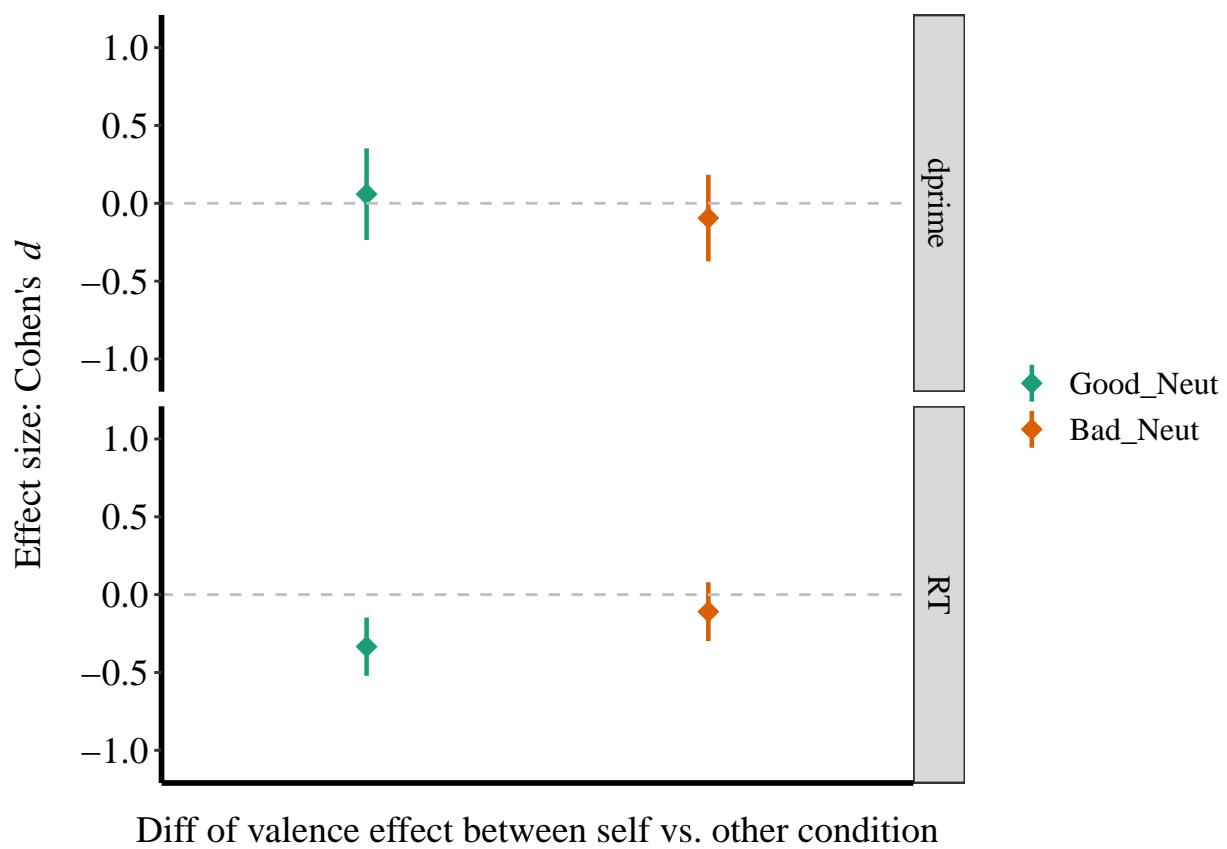


Figure 35. Effect size (Cohen's d) of Valence in Exp4b.

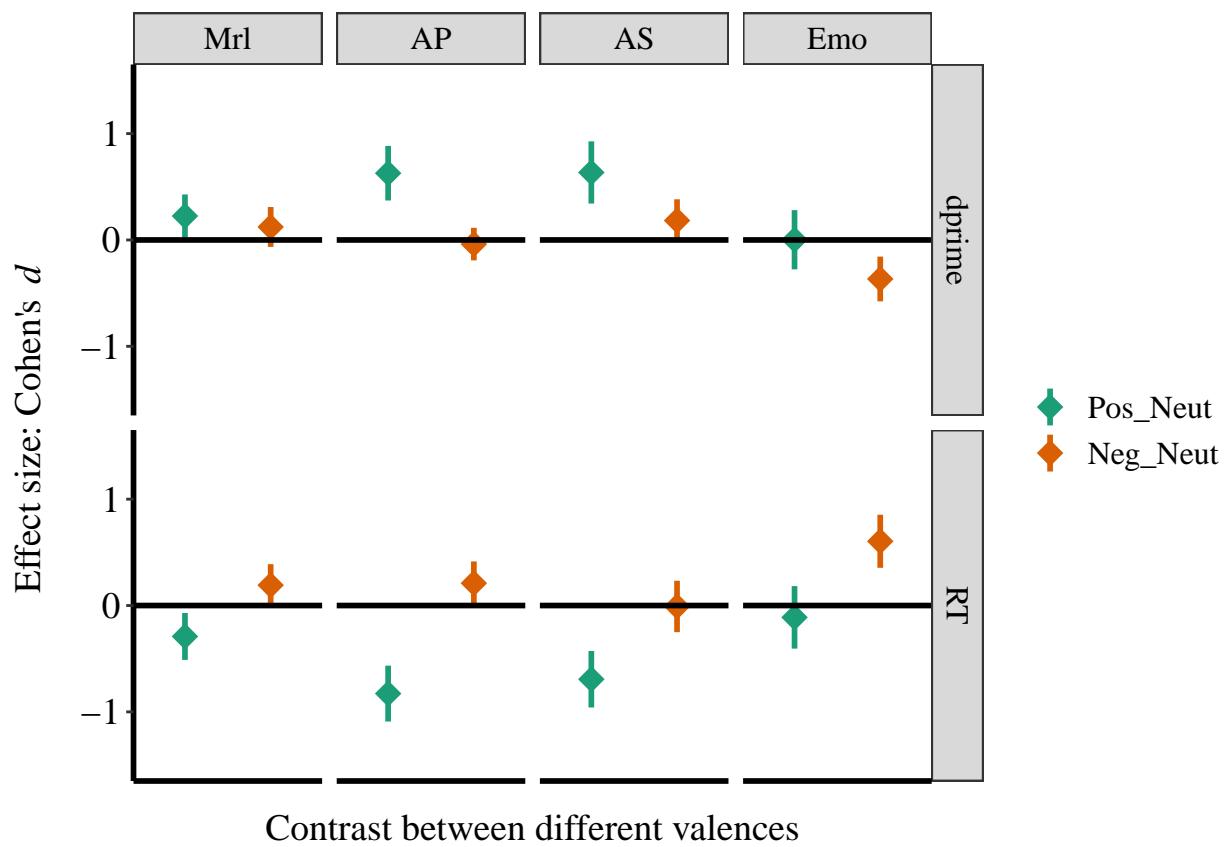


Figure 36. Effect size (Cohen's d) of Valence in Exp5.

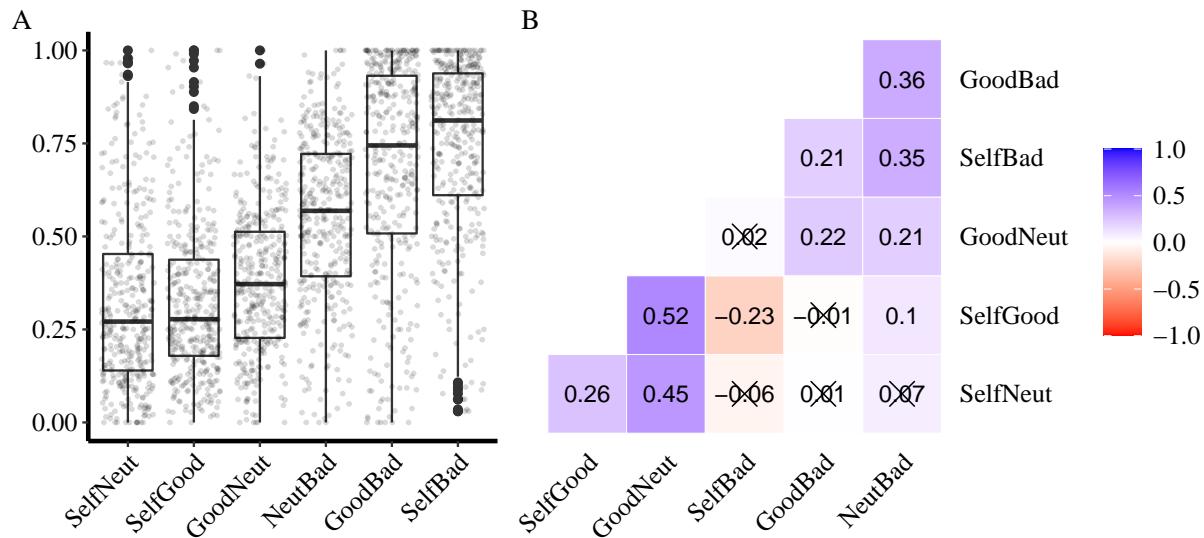


Figure 37. Self-rated personal distance

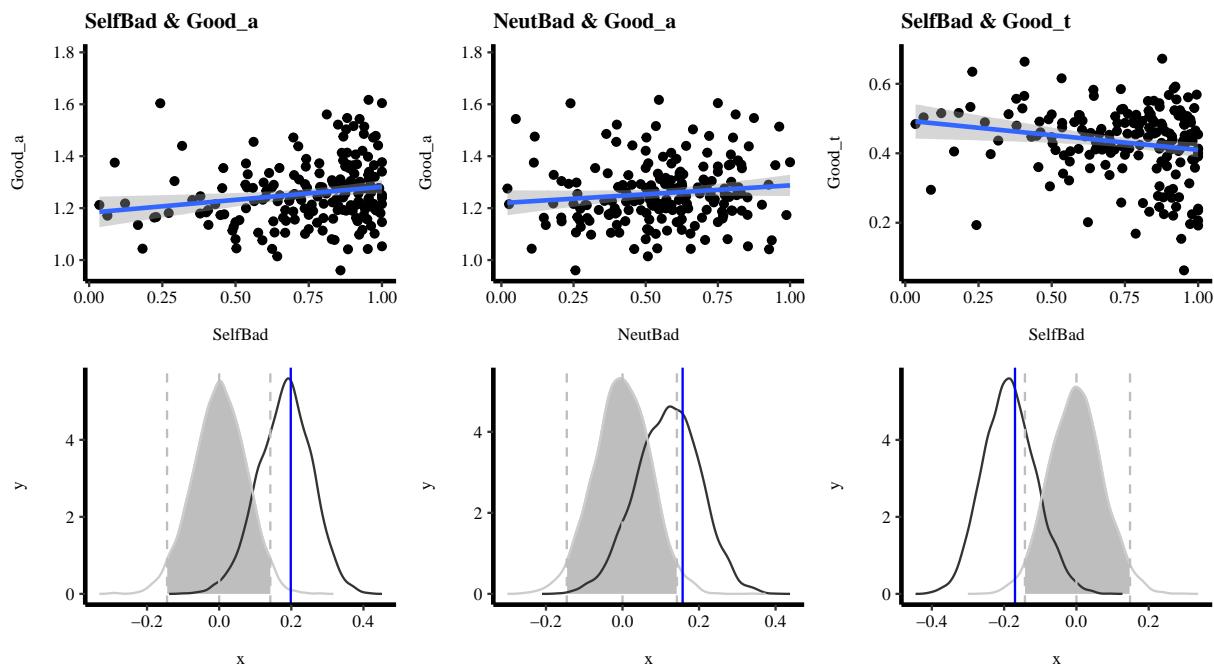


Figure 38. Correlation between moral identity and boundary separation of good condition; moral self-image and drift rate of good condition

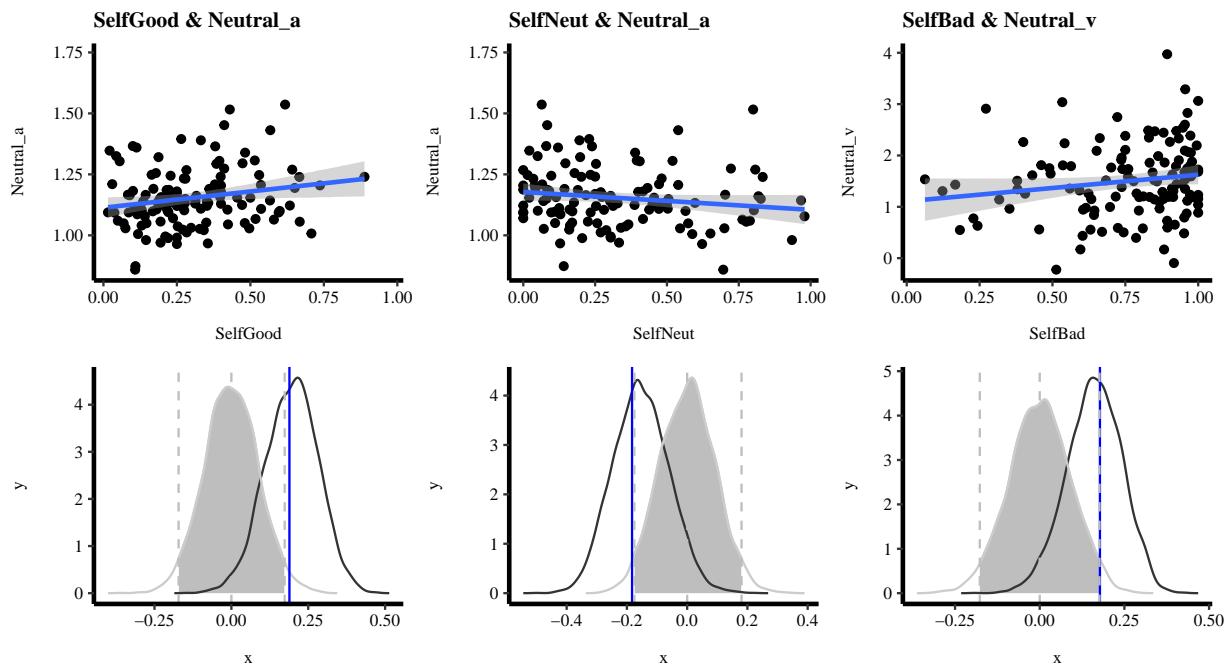


Figure 39. Correlation between personal distance and boundary separation of neutral condition