

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 positive effect was modulated by
25 self-relevance by manipulating the self-referential explicitly and found that this moral
26 positivity effect only occured when the moral valence are self-relevant but evidence to
27 support such effect when the moral valence are other-relevant is weak. We further found
28 that this effect exist even when the self-relevance or the moral valence were presented as a
29 task-irrelevant information, though the effect size become much smaller. We also tested
30 whether the positivity effect only exist in moral domain and found that this effect was not
31 limited to moral domain. Exploratory analyses on task-questionnaire relationship found
32 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 XXXX In perceptual matching, same is faster than different (Farell, 1985; Krueger,

42 1978). Automatic processing (Spruyt & Houwer, 2017)

43 Van Zandt, Colonius, and Proctor (2000): A comparison of two response time models

44 applied to perceptual matching

45 Yakushijin, ReikoJacobs, Robert A (2020), Are People Successful at Learning

46 Sequential Decisions on a Perceptual Matching Task?

47 Schooler, L. J., Shiffrin, R. M., & Raaijmakers, J. G. W. (2001). A Bayesian model

48 for implicit effects in perceptual identification. Psychological Review, 108(1), 257–272.

49 <https://doi.org/10.1037/0033-295X.108.1.257>

50 We reported results from eleven experiments. In first set of experiments, we found

51 that shapes associated with morally positive person label were responded faster and more

52 accurately. In the second set of experiments, we explore the potential role of good self in

53 perceptual matching task and added one more independent variable, we found that the

54 effect was mainly on good self. In the third part we tested whether the morality will

55 automatically binds with person-relevance. Finally, we explore the correlation between

56 behavioral task and questionnaire scores.

57 **Disclosures**

58 We reported all the measurements, analyses, and results in all the experiments in the

59 current study. Participants whose overall accuracy lower than 60% were excluded from

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

61 from the analysis. To have a better overview of the effect reported in this series

62 experiment, we reported the synthesized results in the main text and individual experiment
63 in supplementary materials.

64 All the experiments reported were not pre-registered. Most experiments (1a ~ 6b,
65 except experiment 3b) reported in the current study were first finished between 2014 to
66 2016 in Tsinghua University, Beijing, China. Participants in these experiments were
67 recruited in the local community. To increase the sample size of experiments to 50 or more
68 (Simmons, Nelson, & Simonsohn, 2013), we recruited additional participants in Wenzhou
69 University, Wenzhou, China in 2017 for experiment 1a, 1b, 4a, and 4b. Experiment 3b was
70 finished in Wenzhou University in 2017. To have a better estimation of the effect size, we
71 included the data from two experiments (experiment 7a, 7b) that were reported in Hu,
72 Lan, Macrae, and Sui (2020) (See Table S1 for overview of these experiments).

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

75 **General methods**

76 **Design and Procedure**

77 This series of experiments started to test the effect of instantly acquired true self
78 (moral self) on perceptual decision-making. For this purpose, we used the social associative
79 learning paradigm (or tagging paradigm)(Sui, He, & Humphreys, 2012), in which
80 participants first learned the associations between geometric shapes and labels of person
81 with different moral character (e.g., in first three studies, the triangle, square, and circle
82 and good person, neutral person, and bad person, respectively). The associations of the
83 shapes and label were counterbalanced across participants. After remembered the
84 associations, participants finished a practice phase to familiar with the task, in which they
85 viewed one of the shapes upon the fixation while one of the labels below the fixation and
86 judged whether the shape and the label matched the association they learned. When

87 participants reached 60% or higher accuracy at the end of the practicing session, they
88 started the experimental task which was the same as in the practice phase.

89 The experiment 1a, 1b, 1c, 2, and 6a shared a 2 (matching: match vs. nonmatch) by
90 3 (moral valence: good vs. neutral vs. bad) within-subject design. Experiment 1a was the
91 first one of the whole series studies and 1b, 1c, and 2 were conducted to exclude the
92 potential confounding factors. More specifically, experiment 1b used different Chinese
93 words as label to test whether the effect only occurred with certain familiar words.
94 Experiment 1c manipulated the moral valence indirectly: participants first learned to
95 associate different moral behaviors with different neutral names, after remembered the
96 association, they then performed the perceptual matching task by associating names with
97 different shapes. Experiment 2 further tested whether the way we presented the stimuli
98 influence the effect of valence, by sequentially presenting labels and shapes. Note that part
99 of participants of experiment 2 were from experiment 1a because we originally planned a
100 cross task comparison. Experiment 6a, which shared the same design as experiment 2, was
101 an EEG experiment which aimed at exploring the neural correlates of the effect. But we
102 will focus on the behavioral results of experiment 6a in the current manuscript.

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

114 3b was designed to separate the self-referential trials and other-referential trials. That is,
115 participants finished two different blocks: in the self-referential blocks, they only responded
116 to good-self, neutral-self, and bad-self, with half match trials and half non-match trials; for
117 the other-reference blocks, they only responded to good-other, neutral-other, and
118 bad-other. Experiment 7a and 7b were designed to test the cross task robustness of the
119 effect we observed in the aforementioned experiments (see, Hu et al., 2020). The matching
120 task in these two experiments shared the same design with experiment 3a, but only with
121 two moral valence, i.e., good vs. bad. We didn't include the neutral condition in
122 experiment 7a and 7b because we found that the neutral and bad conditions constantly
123 showed non-significant results in experiment 1 ~ 6.

124 Experiment 4a and 4b were design to test the automaticity of the binding between
125 self/other and moral valence. In 4a, we used only two labels (self vs. other) and two shapes
126 (circle, square). To manipulate the moral valence, we added the moral-related words within
127 the shape and instructed participants to ignore the words in the shape during the task. In
128 4b, we reversed the role of self-reference and valence in the task: participant learnt three
129 labels (good-person, neutral-person, and bad-person) and three shapes (circle, square, and
130 triangle), and the words related to identity, "self" or "other", were presented in the shapes.
131 As in 4a, participants were told to ignore the words inside the shape during the task.

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

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

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

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

157 Data analysis

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

162 ***Classic NHST.***

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

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

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

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

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

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

Bayesian hierarchical generalized linear model (GLM).

190 The classic NHST approach may ignore the uncertainty in estimate of the parameters

191 for SDT (Rouder & Lu, 2005), and using mean RT assumes normal distribution of RT

192 data, which is always not true because RTs distribution is skewed (Rousselet & Wilcox,

193 2019). To better estimate the uncertainty and use a more appropriate model, we also tried

194 Bayesian hierarchical generalized linear model to analyze each experiment's accuracy and

195 RTs data. We used BRMs (Bürkner, 2017) to build the model, which used Stan (Carpenter

196 et al., 2017) to estimate the posterior.

197 In the GLM model, we assume that the accuracy of each trial is Bernoulli distributed

198 (binomial with 1 trial), with probability p_i that $y_i = 1$.

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

199 In the perceptual matching task, the probability p_i can then be modeled as a function of

200 the trial type:

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

201 The outcomes y_i are 0 if the participant responded "nonmatch" on trial i , 1 if they

202 responded "match". The probability of the "match" response for trial i for a participant is

203 p_i . We then write the generalized linear model on the probits (z-scores; Φ , "Phi") of ps . Φ

204 is the cumulative normal density function and maps z scores to probabilities. Given this

205 parameterization, the intercept of the model (β_0) is the standardized false alarm rate

206 (probability of saying 1 when predictor is 0), which we take as our criterion c . The slope of

207 the model (β_1) is the increase of saying 1 when predictor is 1, in z -scores, which is another

208 expression of d' . Therefore, $c = -zHR = -\beta_0$, and $d' = \beta_1$.

209 In each experiment, we had multiple participants, then we need also consider the

210 variations between subjects, i.e., a hierarchical mode in which individual's parameter and

211 the population level parameter are estimated simultaneously. We assume that the

₂₁₂ outcome of each trial is Bernoulli distributed (binomial with 1 trial), with probability p_{ij}
₂₁₃ that $y_{ij} = 1$.

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

₂₁₄ Similarly, the generalized linear model was extended to two levels:

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

₂₁₅ The outcomes y_{ij} are 0 if participant j responded “nonmatch” on trial i , 1 if they
₂₁₆ responded “match”. The probability of the “match” response for trial i for subject j is p_{ij} .
₂₁₇ We again can write the generalized linear model on the probits (z-scores; Φ , “Phi”) of ps .

₂₁₈ The subjective-specific intercepts ($\beta_0 = -zFAR$) and slopes ($\beta_1 = d'$) are described
₂₁₉ by multivariate normal with means and a covariance matrix for the parameters.

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

₂₂₀ For the reaction time, we used the log normal distribution
₂₂₁ ([https://lindeloev.github.io/shiny-rt/#34_\(shifted\)_log-normal](https://lindeloev.github.io/shiny-rt/#34_(shifted)_log-normal)). This distribution has
₂₂₂ two parameters: μ , σ . μ is the mean of the logNormal distribution, and σ is the disperse of
₂₂₃ the distribution. The log normal distribution can be extended to shifted log normal
₂₂₄ distribution, with one more parameter: shift, which is the earliest possible response.

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

₂₂₅ Shifted log-normal distribution:

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

₂₂₆ y_{ij} is the RT of the i th trial of the j th participants.

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

$$\sigma_j \sim Cauchy()$$

227 *Hierarchical drift diffusion model (HDDM).*

To further explore the psychological mechanism under perceptual decision-making, we used HDDM (Wiecki, Sofer, & Frank, 2013) to model our RTs and accuracy data. We used the prior implemented in HDDM, that is, informative priors that constrains parameter estimates to be in the range of plausible values based on past literature (Matzke & 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.

238 *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

²⁴⁹ self-relevance. Cohen's d and its variance were estimated using the following formula
²⁵⁰ (Cooper, Hedges, & Valentine, 2009):

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

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

²⁵¹ M_1 is the mean of the first condition, sd_1 is the standard deviation of the first
²⁵² condition, while M_2 is the mean of the second condition, sd_2 is the standard deviation of
²⁵³ the second condition. r is the correlation coefficient between data from first and second
²⁵⁴ condition. n is the number of data point (in our case the number of participants included
²⁵⁵ in our research).

²⁵⁶ The effect size from each experiment were then synthesized by random effect model
²⁵⁷ using `metafor` (Viechtbauer, 2010). Note that to avoid the cases that some participants
²⁵⁸ participated more than one experiments, we inspected the all available information of
²⁵⁹ participants and only included participants' results from their first participation. As
²⁶⁰ mentioned above, 24 participants were intentionally recruited to participate both exp 1a
²⁶¹ and exp 2, we only included their results from experiment 1a in the meta-analysis.

²⁶² We also estimated the synthesized effect size using Bayesian hierarchical model,
²⁶³ which extended the two-level hierarchical model in each experiment into three-level model,
²⁶⁴ which experiment as an additional level. For SDT, we can use a nested hierarchical model
²⁶⁵ to model all the experiment with similar design:

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

²⁶⁶ where

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

- ²⁶⁷ The outcomes y_{ijk} are 0 if participant j in experiment k responded “nonmatch” on trial i ,
²⁶⁸ 1 if they responded “match”.

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

²⁶⁹ and the experiment level parameter mu_{0k} and mu_{1k} is from a higher order
²⁷⁰ distribution:

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

²⁷¹ in which μ_0 and μ_1 means the population level parameter.

²⁷² This model can be easily expand to three-level model in which participants and
²⁷³ experiments are two group level variable and participants were nested in the experiments.

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

²⁷⁴ y_{ijk} is the RT of the i th trial of the j th participants in the k th experiment.

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

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

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

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

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

Valence effect.

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

Explicit interaction between Valence and self-relevance.

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

Implicit interaction between valence and self-relevance.

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

Specificity of the valence effect.

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

304 ***Behavior-Questionnaire correlation.***

305 Finally, we explored correlation between results from behavioral results and
 306 self-reported measures.

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

312 For the behavioral task part, we derived different indices. First, we used the mean of
 313 the RT and d' from each participants of each condition. Second, we used three parameters
 314 from drift diffusion model: drift rate (v), boundary separation (a), and non
 315 decision-making time (t). Third, we calculated the differences between different conditions
 316 (valence effect: good-self vs. bad-self, good-self vs. neutral-self, bad-self vs. neutral-self;
 317 good-other vs. bad-other, good-other vs. neutral-other, bad-other vs. neutral-other;
 318 Self-reference effect: good-self vs. good-other, neutral-self vs. neutral-other, bad-self
 319 vs. bad-other), as indexed by Cohen's d and standard error (SE) of Cohen's d .

$$Cohen's d_z = \frac{(M_1 - M_2)}{\sqrt{(SD_1^2 + SD_2^2)/2}}$$

320 Given that the task difficulty were different across experiments, we z-transformed all these
 321 indices so that they become unit-free.

322 We used the mean of parameter posterior distribution as the estimate of each
 323 parameter for each participants in the correlation analysis.

324 We used Pearson correlation to quantify the correlation. For those correlation that is
325 significant ($p < 0.05$), we further tested the robustness of the correlation using bootstrap
326 by BootES package (Kirby & Gerlanc, 2013). To avoid false positive, we further determined
327 the threshold for significant by permutation. More specifically, for each pairs that initially
328 with $p < .05$, we randomly shuffle the participants data of each score and calculated the
329 correlation between the shuffled vectors. After repeating this procedure for 5000 times, we
330 choose arrange these 5000 correlation coefficients and use the 95% percentile number as our
331 threshold.

332 **Part 1: Moral valence effect**

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

336 **Experiment 1a**

337 **Methods.**

338 ***Participants.***

339 57 college students (38 female, age = 20.75 ± 2.54 years) participated. 39 of them
340 were recruited from Tsinghua University community in 2014; 18 were recruited from
341 Wenzhou University in 2017. All participants were right-handed except one, and all had
342 normal or corrected-to-normal vision. Informed consent was obtained from all participants
343 prior to the experiment according to procedures approved by the local ethics committees. 6
344 participant's data were excluded from analysis because nearly random level of accuracy,
345 leaving 51 participants (34 female, age = 20.72 ± 2.44 years).

346 ***Stimuli and Tasks.***

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

348 These shapes were paired with three labels (bad person, good person or neutral person).

349 The pairs were counterbalanced across participants.

350 ***Procedure.***

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

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

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

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

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

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

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

358 shape-label pairs, which were subsequently presented, were correct. Each trial started with

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

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

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

362 stage, or it could be a recombination of a shape with a different label, with the shape-label

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

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

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

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

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

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

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

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

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

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

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

374 neutral-person nonmatch). For the participants from Wenzhou University, they finished 6
375 blocks of 120 trials, therefore, 120 trials for each condition.

376 ***Data analysis.***

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

380 **Results.**

381 ***Classic NHST.***

382 *d prime.*

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

385 We found the effect of Valence ($F(1.96, 97.84) = 6.19$, $MSE = 0.27$, $p = .003$,
386 $\hat{\eta}_G^2 = .020$). The post-hoc comparison with multiple comparison correction revealed that
387 the shapes associated with Good-person (2.11, SE = 0.14) has greater *d* prime than shapes
388 associated with Bad-person (1.75, SE = 0.14), $t(50) = 3.304$, $p = 0.0049$. The Good-person
389 condition was also greater than the Neutral-person condition (1.95, SE = 0.16), but didn't
390 reach statistical significant, $t(50) = 1.54$, $p = 0.28$. Neither the Neutral-person condition is
391 significantly greater than the Bad-person condition, $t(50) = 2.109$, $p = .098$.

392 *Reaction times.*

393 We conducted 2 (Matchness: match v. nonmatch) by 3 (Valence: good, neutral, bad)
394 repeated measure ANOVA. We found the main effect of Matchness ($F(1, 50) = 232.39$,
395 $MSE = 948.92$, $p < .001$, $\hat{\eta}_G^2 = .104$), main effect of valence ($F(1.87, 93.31) = 9.62$,
396 $MSE = 1,673.86$, $p < .001$, $\hat{\eta}_G^2 = .016$), and interaction between Matchness and Valence
397 ($F(1.73, 86.65) = 8.52$, $MSE = 1,441.75$, $p = .001$, $\hat{\eta}_G^2 = .011$).

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

399 matched trials, we found the effect of valence . We further examined the effect of valence
 400 for both self and other for matched trials. We found that shapes associated with Good
 401 Person (684 ms, SE = 11.5) responded faster than Neutral (709 ms, SE = 11.5), $t(50) =$
 402 -2.265, $p = 0.0702$ and Bad Person (728 ms, SE = 11.7), $t(50) = -4.41$, $p = 0.0002$), and
 403 the Neutral condition was faster than the Bad condition, $t(50) = -2.495$, $p = 0.0415$). For
 404 non-matched trials, there was no significant effect of Valence ()�.

405 ***Bayesian hierarchical GLM.***

406 *d prime.*

407 We fitted a Bayesian hierarchical GLM for signal detection theory approach. The
 408 results showed that when the shapes were tagged with labels with different moral valence,
 409 the sensitivity (d') and criteria (c) were both influence. For the d' , we found that the
 410 shapes tagged with morally good person (2.46, 95% CI[2.21 2.72]) is greater than shapes
 411 tagged with moral bad (2.07, 95% CI[1.83 2.32]), $P_{PosteriorComparison} = 1$. Shape tagged
 412 with morally good person is also greater than shapes tagged with neutral person (2.23,
 413 95% CI[1.95 2.49]), $P_{PosteriorComparison} = 0.97$. Also, the shapes tagged with neutral
 414 person is greater than shapes tagged with morally bad person, $P_{PosteriorComparison} = 0.92$.

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

420 *Reaction times.*

421 We fitted a Bayesian hierarchical GLM for RTs, with a log-normal distribution as the
 422 link function. We used the posterior distribution of the regression coefficient to make
 423 statistical inferences. As in previous studies, the matched conditions are much faster than
 424 the mismatched trials ($P_{PosteriorComparison} = 1$). We focused on matched trials only, and

425 compared different conditions: Good is faster than the neutral, $P_{PosteriorComparison} = .99$,
426 it was also faster than the Bad condition, $P_{PosteriorComparison} = 1$. And the neutral
427 condition is faster than the bad condition, $P_{PosteriorComparison} = .99$. However, the
428 mismatched trials are largely overlapped. See Figure 2.

429 **HDDM.**

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

437 **Experiment 1b**

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

440 **Method.**

441 **Participants.**

442 72 college students (49 female, age = 20.17 ± 2.08 years) participated. 39 of them
443 were recruited from Tsinghua University community in 2014; 33 were recruited from
444 Wenzhou University in 2017. All participants were right-handed except one, and all had
445 normal or corrected-to-normal vision. Informed consent was obtained from all participants
446 prior to the experiment according to procedures approved by the local ethics committees.
447 20 participant's data were excluded from analysis because nearly random level of accuracy,
448 leaving 52 participants (36 female, age = 20.25 ± 2.31 years).

⁴⁴⁹ **Stimuli and Tasks.** Three geometric shapes (triangle, square, and circle, with 3.7°

⁴⁵⁰ $\times 3.7^\circ$ of visual angle) were presented above a white fixation cross subtending $0.8^\circ \times 0.8^\circ$

⁴⁵¹ of visual angle at the center of the screen. The three shapes were randomly assigned to

⁴⁵² three labels with different moral valence: a morally bad person (" ", ERen), a morally

⁴⁵³ good person (" ", ShanRen) or a morally neutral person (" ", ChangRen). The order of

⁴⁵⁴ the associations between shapes and labels was counterbalanced across participants. Three

⁴⁵⁵ labels used in this experiment is selected based on the rating results from an independent

⁴⁵⁶ survey, in which participants rated the familiarity, frequency, and concreteness of eight

⁴⁵⁷ different words online. Of the eight words, three of them are morally positive (HaoRen,

⁴⁵⁸ ShanRen, Junzi), two of them are morally neutral (ChangRen, FanRen), and three of them

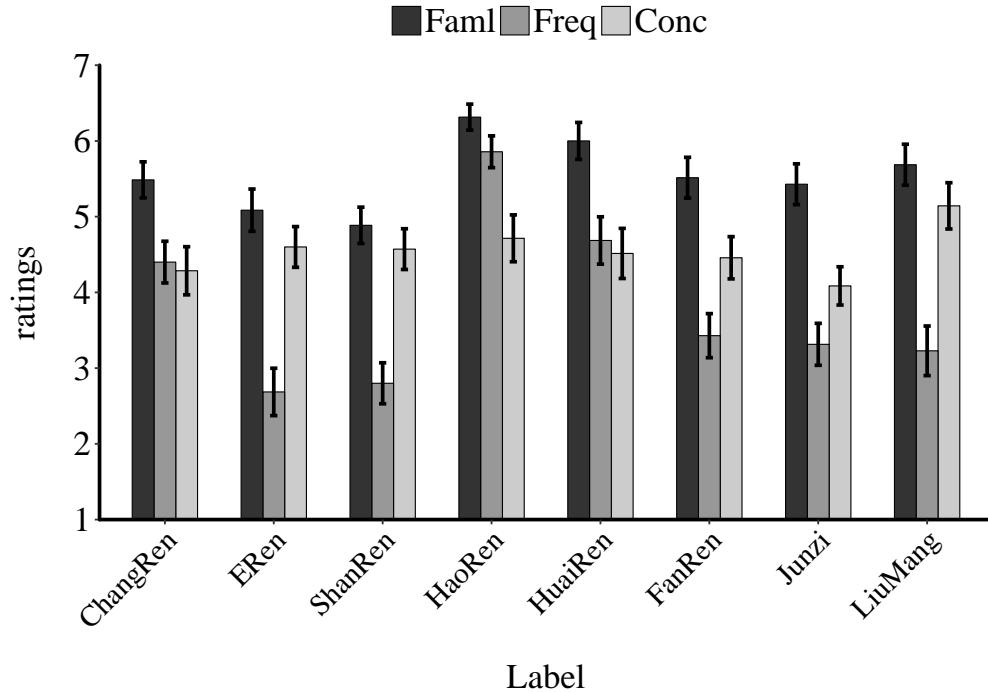
⁴⁵⁹ are morally negative (HuaiRen, ERen, LiuMang). An independent sample consist of 35

⁴⁶⁰ participants (22 females, age 20.6 ± 3.11) were recruited to rate these words. Based on the

⁴⁶¹ ratings (see supplementary materials Figure S1), we selected ShanRen, ChangRen, and

⁴⁶² ERen to represent morally positive, neutral, and negative person.

Ratings for each label



Procedure.

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

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

468 **Results.**

469 **NHST.**

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

471 *d* prime.

472 Repeated measures ANOVA revealed main effect of valence, $F(1.83, 93.20) = 14.98$,
 473 $MSE = 0.18$, $p < .001$, $\hat{\eta}_G^2 = .053$. Paired t test showed that the Good-Person condition
 474 (1.87 ± 0.102) was with greater *d* prime than Neutral condition (1.44 ± 0.101 , $t(51) =$
 475 5.945 , $p < 0.001$). We also found that the Bad-Person condition (1.67 ± 0.11) has also
 476 greater *d* prime than neutral condition , $t(51) = 3.132$, $p = 0.008$). There Good-person
 477 condition was also slightly greater than the bad condition, $t(51) = 2.265$, $p = 0.0701$.

478 *Reaction time.*

479 We found interaction between Matchness and Valence ($F(1.95, 99.31) = 19.71$,
 480 $MSE = 960.92$, $p < .001$, $\hat{\eta}_G^2 = .031$) and then analyzed the matched trials and
 481 mismatched trials separately, as in experiment 1a. For matched trials, we found the effect
 482 of valence $F(1.94, 99.10) = 33.97$, $MSE = 1,343.19$, $p < .001$, $\hat{\eta}_G^2 = .115$. Post-hoc *t*-tests
 483 revealed that shapes associated with Good Person (684 ± 8.77) were responded faster than
 484 Neutral-Person (740 ± 9.84), ($t(51) = -8.167$, $p < 0.001$) and Bad Person (728 ± 9.15),
 485 $t(51) = -5.724$, $p < 0.0001$). While there was no significant differences between Neutral and
 486 Bad-Person condition ($t(51) = 1.686$, $p = 0.221$). For non-matched trials, there was no
 487 significant effect of Valence ($F(1.90, 97.13) = 1.80$, $MSE = 430.15$, $p = .173$, $\hat{\eta}_G^2 = .003$).

488 **BGLM.**

489 *Signal detection theory analysis of accuracy.*

490 We fitted a Bayesian hierarchical GLM for SDT. The results showed that when the
491 shapes were tagged with labels with different moral valence, the sensitivity (d') and criteria
492 (c) were both influence. For the d' , we found that the shapes tagged with morally good
493 person (2.46, 95% CI[2.21 2.72]) is greater than shapes tagged with moral bad (2.07, 95%
494 CI[1.83 2.32]), $P_{PosteriorComparison} = 1$. Shape tagged with morally good person is also
495 greater than shapes tagged with neutral person (2.23, 95% CI[1.95 2.49]),
496 $P_{PosteriorComparison} = 0.97$. Also, the shapes tagged with neutral person is greater than
497 shapes tagged with morally bad person, $P_{PosteriorComparison} = 0.92$.

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

503 *Reaction time.*

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

512 **HDDM.**

513 We found that the shapes tagged with good person has higher drift rate and higher
514 boundary separation than shapes tagged with both neutral and bad person. Also, the

515 shapes tagged with neutral person has a higher drift rate than shapes tagged with bad
516 person, but not for the boundary separation. Finally, we found that shapes tagged with
517 bad person had longer non-decision time (see figure 6).

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

523 **Experiment 1c**

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

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

528 **Method.**

529 ***Participants.***

530 23 college students (15 female, age = 22.61 ± 2.62 years) participated. All of them
531 were recruited from Tsinghua University community in 2014. Informed consent was
532 obtained from all participants prior to the experiment according to procedures approved by
533 the local ethics committees. No participant was excluded because they overall accuracy
534 were above 0.6.

535 ***Stimuli and Tasks.***

536 Three geometric shapes (triangle, square, and circle, with $3.7^\circ \times 3.7^\circ$ of visual angle)
537 were presented above a white fixation cross subtending $0.8^\circ \times 0.8^\circ$ of visual angle at the
538 center of the screen. The three most common names were chosen, which are neutral in
539 moral valence before the manipulation. Three names (Zhang, Wang, Li) were first paired

540 with three paragraphs of behavioral description. Each description includes one sentence of
541 biographic information and four sentences that describing the moral behavioral under that
542 name. To assess the that these three descriptions represented good, neutral, and bad
543 valence, we collected the ratings of three person on six dimensions: morality, likability,
544 trustworthiness, dominance, competence, and aggressiveness, from an independent sample
545 ($n = 34$, 18 female, age = 19.6 ± 2.05). The rating results showed that the person with
546 morally good behavioral description has higher score on morality ($M = 3.59$, $SD = 0.66$)
547 than neutral ($M = 0.88$, $SD = 1.1$), $t(33) = 12.94$, $p < .001$, and bad conditions ($M = -3.4$,
548 $SD = 1.1$), $t(33) = 30.78$, $p < .001$. Neutral condition was also significant higher than bad
549 conditions $t(33) = 13.9$, $p < .001$ (See supplementary materials).

550 **Procedure.**

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

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

564 **Results.** Figure 7 shows d prime and reaction times of experiment 1c. We
565 conducted same analysis as in Experiment 1a. Our analysis didn't show effect of valence

566 on d prime, $F(1.93, 42.56) = 0.23$, $MSE = 0.41$, $p = .791$, $\hat{\eta}_G^2 = .005$. Neither the effect of
 567 valence on RT ($F(1.63, 35.81) = 0.22$, $MSE = 2,212.71$, $p = .761$, $\hat{\eta}_G^2 = .001$) or
 568 interaction between valence and matchness on RT ($F(1.79, 39.43) = 1.20$,
 569 $MSE = 1,973.91$, $p = .308$, $\hat{\eta}_G^2 = .005$).

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

571 We fitted a Bayesian hierarchical GLM for SDT. The results showed that when the
 572 shapes were tagged with labels with different moral valence, the sensitivity (d') and criteria
 573 (c) were both influenced. For the d' , we found that the shapes tagged with morally good
 574 person (2.30, 95% CI[1.93 2.70]) is greater than shapes tagged with moral bad (2.11, 95%
 575 CI[1.83 2.42]), $P_{PosteriorComparison} = 0.8$. Shape tagged with morally good person is also
 576 greater than shapes tagged with neutral person (2.16, 95% CI[1.88 2.45]),
 577 $P_{PosteriorComparison} = 0.75$.

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

583 ***Reaction time.***

584 We fitted a Bayesian hierarchical GLM for RTs, with a log-normal distribution as the
 585 link function. We used the posterior distribution of the regression coefficient to make
 586 statistical inferences. As in previous studies, the matched conditions are much faster than
 587 the mismatched trials ($P_{PosteriorComparison} = .75$). We focused on matched trials only, and
 588 compared different conditions: Good () is not faster than the neutral (),
 589 $P_{PosteriorComparison} = .5$, it was faster than the Bad condition (),
 590 $P_{PosteriorComparison} = .88$. And the neutral condition is faster than the bad condition,
 591 $P_{PosteriorComparison} = .95$. However, the mismatched trials are largely overlapped.

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

599 **Experiment 2: Sequential presenting**

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

612 **Method.**

613 **Participants.**

614 35 participants (17 female, age = 21.66 ± 3.03) were recruited. 24 of them had
615 participated in Experiment 1a (9 male, mean age = 21.9, s.d. = 2.9), and the time gap
616 between these experiment 1a and experiment 2 is at least six weeks. The results of 1

617 participants were excluded from analysis because of less than 60% overall accuracy,
618 remains 34 participants (17 female, age = 21.74 ± 3.04).

619 ***Procedure.***

620 In Experiment 2, the sequential presenting makes the matching task much easier than
621 experiment 1. To avoid ceiling effect on behavioral data, we did a few pilot experiments to
622 get optimal parameters, i.e., the conditions under which participant have similar accuracy
623 as in Experiment 1 (around 70 ~ 80% accuracy). In the final procedure, the label (good
624 person, bad person, or neutral person) was presented for 50 ms and then masked by a
625 scrambled image for 200 ms. A geometric shape followed the scrambled mask for 50 ms in
626 a noisy background (which was produced by first decomposing a square with $\frac{3}{4}$ gray area
627 and $\frac{1}{4}$ white area to small squares with a size of 2×2 pixels and then re-combine these
628 small pieces randomly), instead of pure gray background in Experiment 1. After that, a
629 blank screen was presented 1100 ms, during which participants should press a button to
630 indicate the label and the shape match the original association or not. Feedback was given,
631 as in study 1. The next trial then started after 700 ~ 1100 ms blank. Other aspects of
632 study 2 were identical to study 1.

633 ***Data analysis.***

634 Data was analyzed as in study 1a.

635 **Results.**

636 ***NHST.***

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

639 ***d prime.***

640 There was evidence for the main effect of valence, $F(1.83, 60.36) = 14.41$,
641 $MSE = 0.23$, $p < .001$, $\hat{\eta}_G^2 = .066$. Paired t test showed that the Good-Person condition

642 (2.79 ± 0.17) was with greater d prime than Netural condition (2.21 ± 0.16, $t(33) = 4.723$,
 643 $p = 0.001$) and Bad-person condition (2.41 ± 0.14), $t(33) = 4.067$, $p = 0.008$). There was
 644 no-significant difference between Neutral-person and Bad-person conidition, $t(33) = -1.802$,
 645 $p = 0.185$.

646 *Reaction time.*

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

649 We found interaction between Matchness and Valence ($F(1.99, 65.70) = 9.53$,
 650 $MSE = 605.36$, $p < .001$, $\hat{\eta}_G^2 = .017$) and then analyzed the matched trials and
 651 mismatched trials separately, as in experiment 1a. For matched trials, we found the effect
 652 of valence $F(1.99, 65.76) = 10.57$, $MSE = 1,192.65$, $p < .001$, $\hat{\eta}_G^2 = .067$. Post-hoc t -tests
 653 revealed that shapes associated with Good Person (548 ± 9.4) were responded faster than
 654 Neutral-Person (582 ± 10.9), ($t(33) = -3.95$, $p = 0.0011$) and Bad Person (582 ± 10.2),
 655 $t(33) = -3.9$, $p = 0.0013$). While there was no significant differences between Neutral and
 656 Bad-Person condition ($t(33) = -0.01$, $p = 0.999$). For non-matched trials, there was no
 657 significant effect of Valence ($F(1.99, 65.83) = 0.17$, $MSE = 489.80$, $p = .843$, $\hat{\eta}_G^2 = .001$).

658 **BGLMM.**

659 *Signal detection theory analysis of accuracy.*

660 We fitted a Bayesian hierarchical GLM for SDT. The results showed that when the
 661 shapes were tagged with labels with different moral valence, the sensitivity (d') and criteria
 662 (c) were both influence. For the d' , we found that the shapes tagged with morally good
 663 person (2.46, 95% CI[2.21 2.72]) is greater than shapes tagged with moral bad (2.07, 95%
 664 CI[1.83 2.32]), $P_{PosteriorComparison} = 1$. Shape tagged with morally good person is also
 665 greater than shapes tagged with neutral person (2.23, 95% CI[1.95 2.49]),
 666 $P_{PosteriorComparison} = 0.97$. Also, the shapes tagged with neutral person is greater than
 667 shapes tagged with morally bad person, $P_{PosteriorComparison} = 0.92$.

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

673 *Reaction times.*

674 We fitted a Bayesian hierarchical GLM for RTs, with a log-normal distribution as the
 675 link function. We used the posterior distribution of the regression coefficient to make
 676 statistical inferences. As in previous studies, the matched conditions are much faster than
 677 the mismatched trials ($P_{PosteriorComparison} = .75$). We focused on matched trials only, and
 678 compared different conditions: Good () is not faster than the neutral (),
 679 $P_{PosteriorComparison} = .5$, it was faster than the Bad condition (),
 680 $P_{PosteriorComparison} = .88$. And the neutral condition is faster than the bad condition,
 681 $P_{PosteriorComparison} = .95$. However, the mismatched trials are largely overlapped.

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

690 **Discussion**

691 In this experiment, we repeated the results pattern that the positive moral valenced
 692 stimuli has an advantage over the neutral or the negative valence association. Moreover,

693 with a cross-task analysis, we didn't found evidence that the experiment task interacted
694 with moral valence, suggesting that the effect might not be effect by experiment task.
695 These findings suggested that the facilitation effect of positive moral valence is robust and
696 not affected by task. This robust effect detected by the associative learning is unexpected.

Results

698 Effect of moral valence

699 In this part, we synthesized results from experiment 1a, 1b, 1c, 2, 5 and 6a. Data
700 from 192 participants were included in these analyses. We found differences between
701 positive and negative conditions on RT was Cohen's $d = -0.58 \pm 0.06$, 95% CI [-0.70 -0.47];
702 on d' was Cohen's $d = 0.24 \pm 0.05$, 95% CI [0.15 0.34]. The effect was also observed
703 between positive and neutral condition, RT: Cohen's $d = -0.44 \pm 0.10$, 95% CI [-0.63
704 -0.25]; d' : Cohen's $d = 0.31 \pm 0.07$, 95% CI [0.16 0.45]. And the difference between neutral
705 and bad conditions are not significant, RT: Cohen's $d = 0.15 \pm 0.07$, 95% CI [0.00 0.30];
706 d' : Cohen's $d = 0.07 \pm 0.07$, 95% CI [-0.08 0.21]. See Figure 13 left panel.

707 Interaction between valence and self-reference

708 In this part, we combined the experiments that explicitly manipulated the
709 self-reference and valence, which includes 3a, 3b, 6b, 7a, and 7b. For the positive versus
710 negative contrast, data were from five experiments with 178 participants; for positive
711 versus neutral and neutral versus negative contrasts, data were from three experiments (712 3a, 3b, and 6b) with 108 participants.

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

For the self-referential condition, we found the same pattern as in the first part of results. That is we found significant differences between positive and neutral as well as positive and negative, but not neutral and negative. The effect size of RT between positive and negative is Cohen's $d = -0.89 \pm 0.12$, 95% CI [-1.11 -0.66]; on d' was Cohen's $d = 0.61 \pm 0.09$, 95% CI [0.44 0.78]. The effect was also observed between positive and neutral condition, RT: Cohen's $d = -0.76 \pm 0.13$, 95% CI [-1.01 -0.50]; d' : Cohen's $d = 0.69 \pm 0.14$, 95% CI [0.42 0.96]. And the difference between neutral and bad conditions are not significant, RT: Cohen's $d = 0.03 \pm 0.13$, 95% CI [-0.22 0.29]; d' : Cohen's $d = 0.08 \pm 0.08$, 95% CI [-0.07 0.24]. See Figure 13 the middle panel.

For the other-referential condition, we found that only the difference between positive and negative on RT was significant, all the other conditions were not. The effect size of RT between positive and negative is Cohen's $d = -0.28 \pm 0.05$, 95% CI [-0.38 -0.17]; on d' was 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 13 right panel.

Generalizability of the valence effect

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

For experiment 4a, when self-reference was the target and moral valence was task-irrelevant, we found that only under the implicit self-referential condition, i.e., when the moral words were presented as task irrelevant stimuli, there was the main effect of valence and interaction between valence and reference for both d prime and RT (See supplementary results for the detailed statistics). For d prime, we found good-self

742 condition (2.55 ± 0.86) had higher d prime than bad-self condition (2.38 ± 0.80); good self
743 condition was also higher than neutral self (2.45 ± 0.78) but there was not statistically
744 significant, while the neutral-self condition was higher than bad self condition and not
745 significant neither. For reaction times, good-self condition (654.26 ± 67.09) were faster
746 relative to bad-self condition (665.64 ± 64.59), and over neutral-self condition ($664.26 \pm$
747 64.71). The difference between neutral-self and bad-self conditions were not significant.
748 However, for the other-referential condition, there was no significant differences between
749 different valence conditions. See Figure 14.

750 For experiment 4b, when valence was the target and the identity was task-irrelevant,
751 we found a strong valence effect (see supplementary results and Figure 15, Figure 16).

752 In this experiment, the advantage of good-self condition can only be disentangled by
753 comparing the self-referential and other-referential conditions. Therefore, we calculated the
754 differences between the valence effect under self-referential and other referential conditions
755 and used the weighted variance as the variance of this differences. We found this
756 modulation effect on RT. The valence effect of RT was stronger in self-referential than
757 other-referential for the Good vs. Neutral condition (-0.33 ± 0.01), and to a less extent the
758 Good vs. Bad condition (-0.17 ± 0.01). While the size of the other effect's CI included
759 zero, suggestion those effects didn't differ from zero. See Figure 17.

760 Specificity of valence effect

761 In this part, we analyzed the results from experiment 5, which included positive,
762 neutral, and negative valence from four different domains: morality, emotion, aesthetics of
763 human, and aesthetics of scene. We found interaction between valence and domain for both
764 d prime and RT (match trials). A common pattern appeared in all four domains: each
765 domain showed a binary results instead of gradient on both d prime and RT. For morality,
766 aesthetics of human, and aesthetics of scene, the positive conditions had advantages over

767 both neutral and negative conditions (greater d' prime and faster RT), and neutral and
768 negative conditions didn't differ from each other. But for the emotional stimuli, it was the
769 positive and neutral had advantage over negative conditions, while positive and neutral
770 conditions were not significantly different. See supplementary materials for detailed
771 statistics. Also note that the effect size in moral domain is smaller than the aesthetic
772 domains (beauty of people and beauty of scene). See Figure 18.

773 **Self-reported personal distance**

774 See Figure 19.

775 **Correlation analyses**

776 The reliability of questionnaires can be found in (Liu et al., 2020). We calculated the
777 correlation between the data from behavioral task and the questionnaire data.

778 We focused on the task-questionnaire correlation, the results revealed that the score
779 from three questionnaire are related to behavioral responses data. First, the external moral
780 identity is positively correlated with boundary separation of moral good condition,
781 $r = 0.194$, 95% CI [0.023 0.350]); the moral self image is positively correlated with the drift
782 rate ($r = 0.191$, 95% CI [-0.016 0.354]) of the morally good condition. See Figure 20.

783 Second, we found the personal distance between self and good is positively correlated
784 with the boundary separation of neutral condition and the self-neutral distance is
785 negatively correlated with the boundary separation of neutral condition. See figure 21

786 Third, we found the self esteem score was negative correlated with the d' of bad
787 conditions ($r = -0.16$, 95% CI [-0.277 -0.038]) and the neutral conditions ($r = -.197$, 95%
788 CI [-0.348 -0.026]). See Figure 22.

789 We also explored the correlation between behavioral data and questionnaire scores
790 separately for experiments with and without self-referential. For experiments without

791 self-referential (Valence effect), we found the personal distance between Good-person and
792 self is positively correlated with boundary separation of good conditions, $r = 0.292$, 95%
793 [0.071 0.485]. also personal distance between the bad and neutral person is positively
794 correlated with non-responding time of bad and neutral conditions, $r = 0.249$, 0.233,
795 respectively.

796 For experiments with self-referential (Valence effect for the self), we found self-esteem
797 is negatively correlated with d prime of neutral condition, $r = -0.272$, [-0.468 -0.052], the
798 self-good distance is positively correlated with d prime for Bad condition, $r = 0.185$,
799 95%CI[0.004 0.354].

800 Discussion

801 References

- 802 Anderson, E., Siegel, E. H., Bliss-Moreau, E., & Barrett, L. F. (2011). The visual impact
803 of gossip. *Science*, 332(6036), 1446–1448. <https://doi.org/10.1126/science.1201574>
- 804 Brainard, D. H. (1997). The psychophysics toolbox. *Spatial Vision*, 10(4), 433–436.
805 Journal Article.
- 806 Bürkner, P.-C. (2017). Brms: An r package for bayesian multilevel models using stan.
807 *Journal of Statistical Software; Vol 1, Issue 1 (2017)*. Journal Article. Retrieved
808 from
809 <https://www.jstatsoft.org/v080/i01> <http://dx.doi.org/10.18637/jss.v080.i01>
- 810 Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M., ...
811 Riddell, A. (2017). Stan: A probabilistic programming language. *Journal of
812 Statistical Software*, 76(1). Journal Article. <https://doi.org/10.18637/jss.v076.i01>
- 813 Cooper, H., Hedges, L. V., & Valentine, J. C. (2009). *The handbook of research synthesis
814 and meta-analysis* (2nd ed.). Book, New York: Sage.

- 815 DeCarlo, L. T. (1998). Signal detection theory and generalized linear models. *Psychological
816 Methods*, 3(2), 186–205. Journal Article. <https://doi.org/10.1037/1082-989X.3.2.186>
- 817 Farrell, B. (1985). "Same"—"different" judgments: A review of current controversies in
818 perceptual comparisons. *Psychological Bulletin*, 98(3), 419–456. Journal Article.
819 <https://doi.org/10.1037/0033-2909.98.3.419>
- 820 Gantman, A. P., & Van Bavel, J. J. (2014). The moral pop-out effect: Enhanced
821 perceptual awareness of morally relevant stimuli. *Cognition*, 132(1), 22–29.
822 <https://doi.org/10.1016/j.cognition.2014.02.007>
- 823 Goh, J. X., Hall, J. A., & Rosenthal, R. (2016). Mini meta-analysis of your own studies:
824 Some arguments on why and a primer on how. *Social and Personality Psychology
825 Compass*, 10(10), 535–549. Journal Article. <https://doi.org/10.1111/spc3.12267>
- 826 Hu, C.-P., Lan, Y., Macrae, C. N., & Sui, J. (2020). Good me bad me: Does valence
827 influence self-prioritization during perceptual decision-making? *Collabra: Psychology*,
828 6(1), 20. Journal Article. <https://doi.org/10.1525/collabra.301>
- 829 Kirby, K. N., & Gerlanc, D. (2013). BootES: An r package for bootstrap confidence
830 intervals on effect sizes. *Behavior Research Methods*, 45(4), 905–927.
831 <https://doi.org/10.3758/s13428-013-0330-5>
- 832 Krueger, L. E. (1978). A theory of perceptual matching. *Psychological Review*, 85(4),
833 278–304. Journal Article. <https://doi.org/10.1037/0033-295X.85.4.278>
- 834 Liu, Q., Wang, F., Yan, W., Peng, K., Sui, J., & Hu, C.-P. (2020). Questionnaire data from
835 the revision of a chinese version of free will and determinism plus scale. *Journal of
836 Open Psychology Data*, 8(1), 1. Journal Article. <https://doi.org/10.5334/jopd.49/>
- 837 Matzke, D., & Wagenmakers, E.-J. (2009). Psychological interpretation of the ex-gaussian
838 and shifted wald parameters: A diffusion model analysis. *Psychonomic Bulletin &
839 Review*, 16(5), 798–817. <https://doi.org/10.3758/PBR.16.5.798>

- 840 Pelli, D. G. (1997). The videotoolbox software for visual psychophysics: Transforming
841 numbers into movies. *Spatial Vision*, 10(4), 437–442. Journal Article.
- 842 Rouder, J. N., & Lu, J. (2005). An introduction to bayesian hierarchical models with an
843 application in the theory of signal detection. *Psychonomic Bulletin & Review*,
844 12(4), 573–604. Journal Article. <https://doi.org/10.3758/bf03196750>
- 845 Rousselet, G. A., & Wilcox, R. R. (2019). Reaction times and other skewed distributions:
846 Problems with the mean and the median. *Meta-Psychology*. preprint.
847 <https://doi.org/10.1101/383935>
- 848 Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2013). Life after p-hacking. Conference
849 Proceedings. <https://doi.org/10.2139/ssrn.2205186>
- 850 Spruyt, A., & Houwer, J. D. (2017). On the automaticity of relational stimulus processing:
851 The (extrinsic) relational simon task. *PLoS One*, 12(10), e0186606. Journal Article.
852 <https://doi.org/10.1371/journal.pone.0186606>
- 853 Stanislaw, H., & Todorov, N. (1999). Calculation of signal detection theory measures.
854 *Behavior Research Methods, Instruments, & Computers*, 31(1), 137–149. Journal
855 Article. <https://doi.org/10.3758/BF03207704>
- 856 Sui, J., He, X., & Humphreys, G. W. (2012). Perceptual effects of social salience: Evidence
857 from self-prioritization effects on perceptual matching. *Journal of Experimental
858 Psychology: Human Perception and Performance*, 38(5), 1105–1117. Journal
859 Article. <https://doi.org/10.1037/a0029792>
- 860 Van Zandt, T., Colonius, H., & Proctor, R. W. (2000). A comparison of two response time
861 models applied to perceptual matching. *Psychonomic Bulletin & Review*, 7(2),
862 208–256. <https://doi.org/10.3758/BF03212980>
- 863 Wiecki, T. V., Sofer, I., & Frank, M. J. (2013). HDDM: Hierarchical bayesian estimation of
864 the drift-diffusion model in python. *Frontiers in Neuroinformatics*, 7.

865

<https://doi.org/10.3389/fninf.2013.00014>

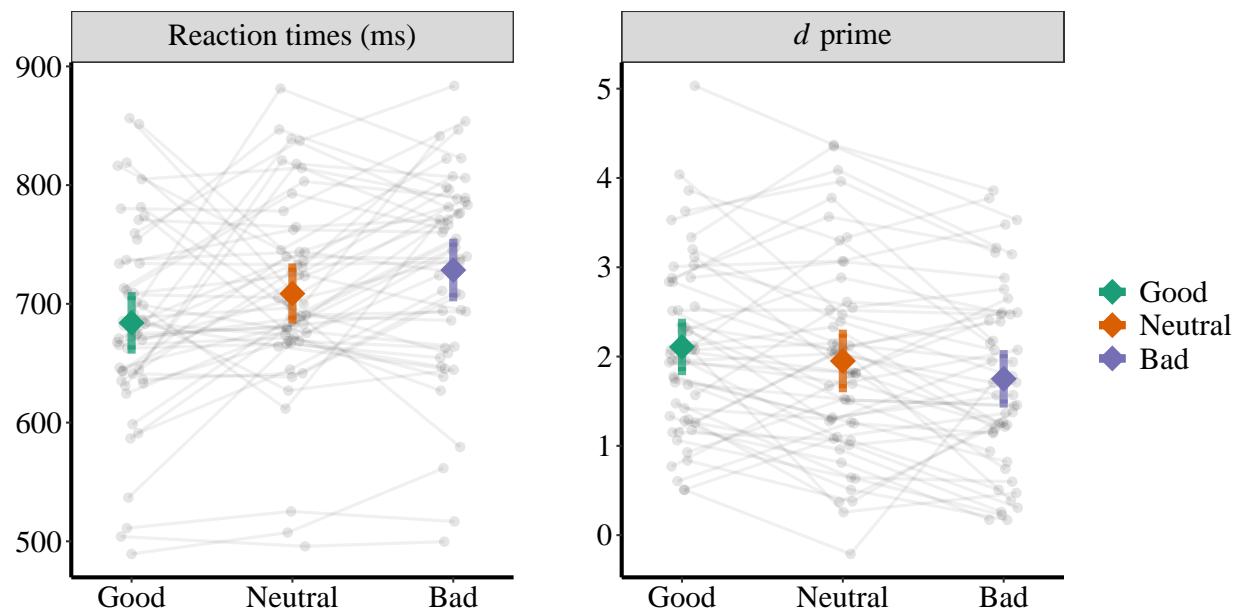


Figure 1. RT and d prime of Experiment 1a.

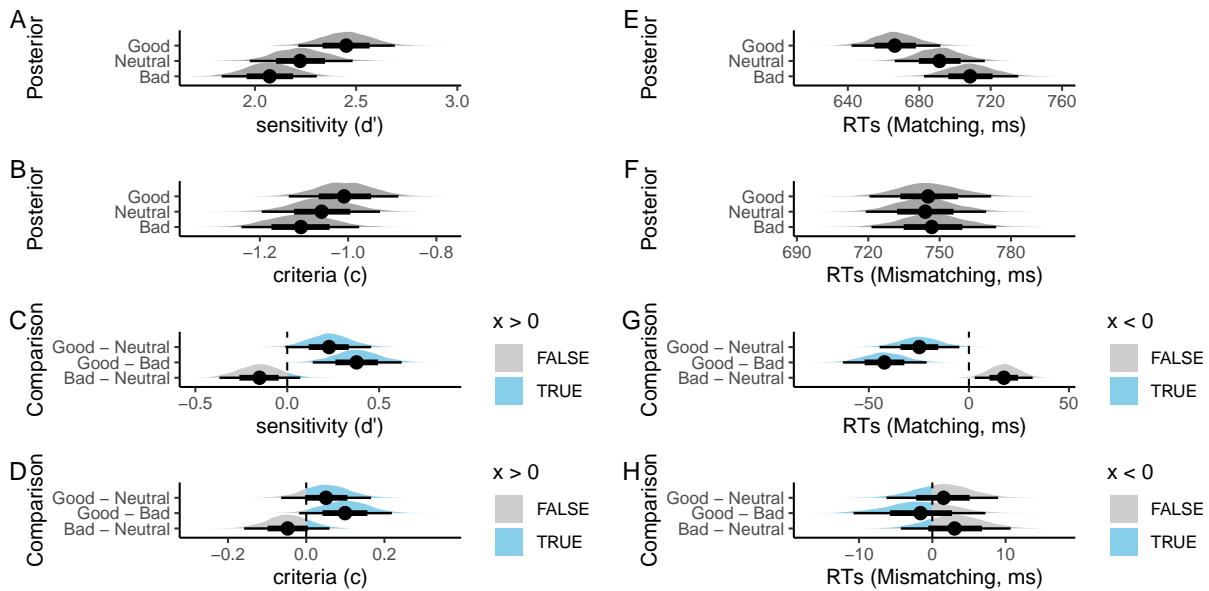


Figure 2. Exp1a: Results of Bayesian GLM analysis.

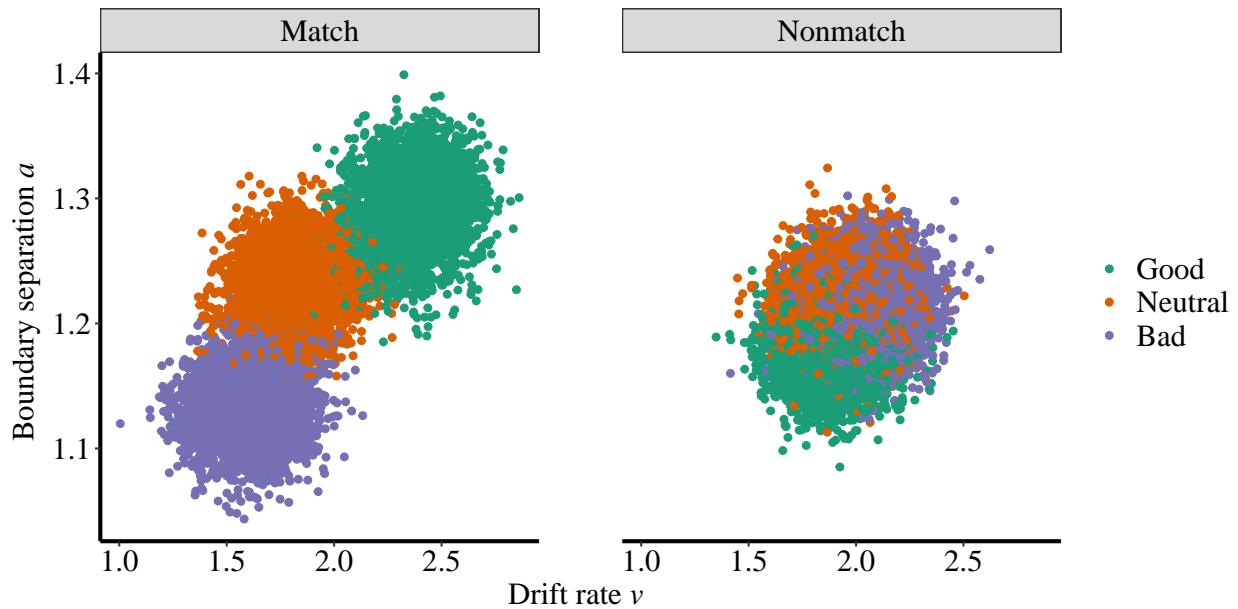


Figure 3. Exp1a: Results of HDDM.

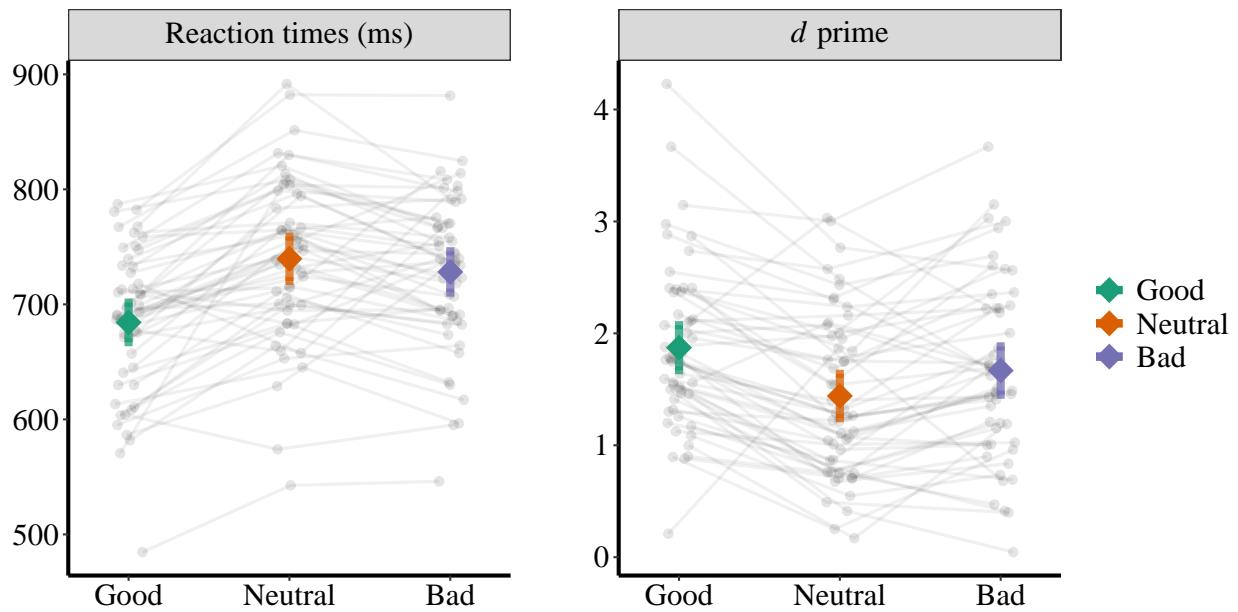


Figure 4. RT and d prime of Experiment 1b.

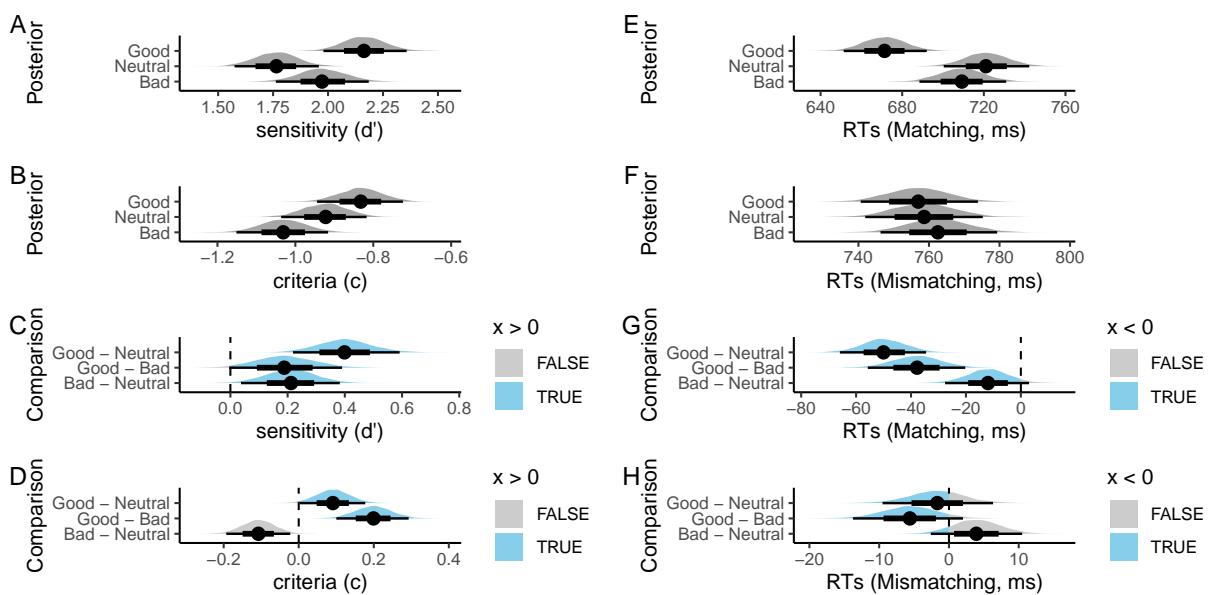


Figure 5. Exp1b: Results of Bayesian GLM analysis.

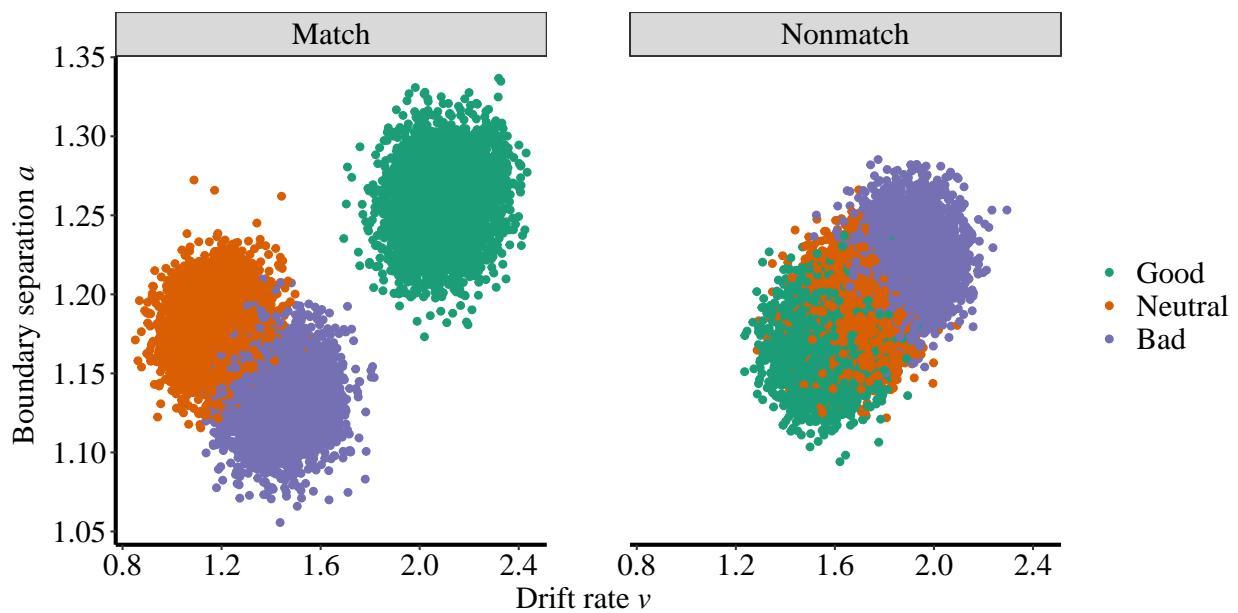


Figure 6. Exp1b: Results of HDDM.

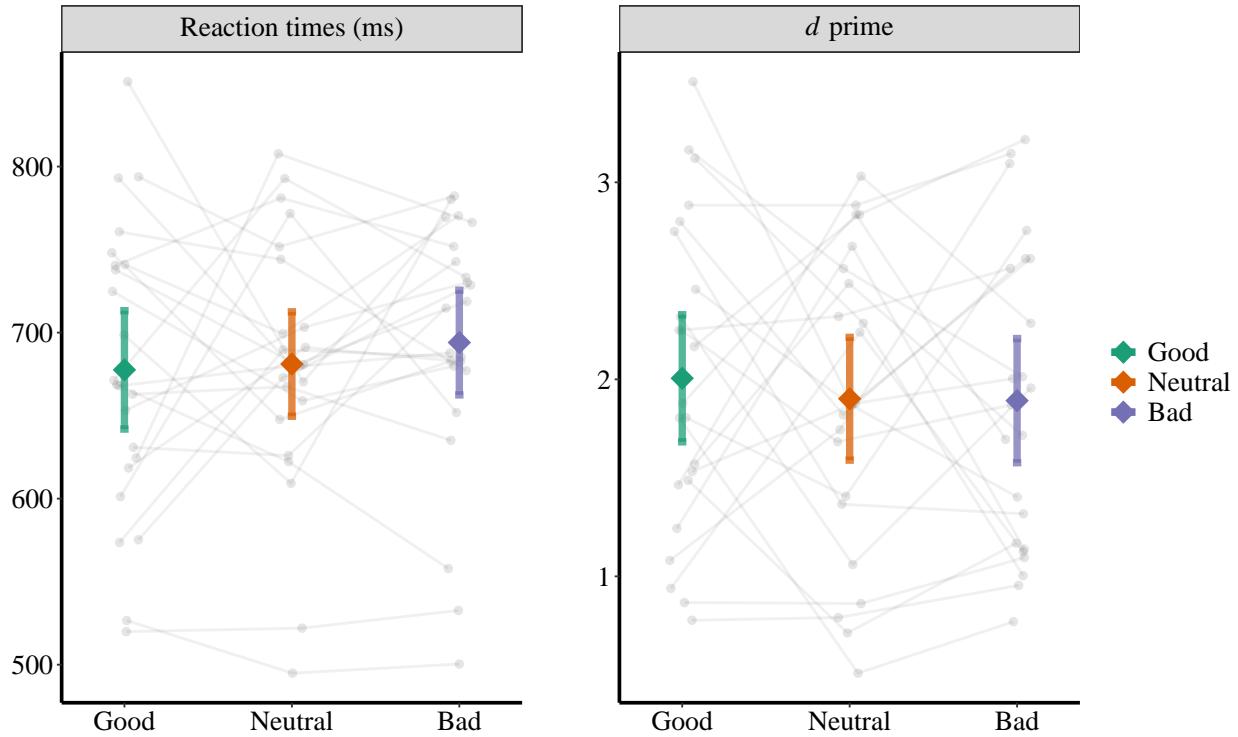


Figure 7. RT and d' prime of Experiment 1c.

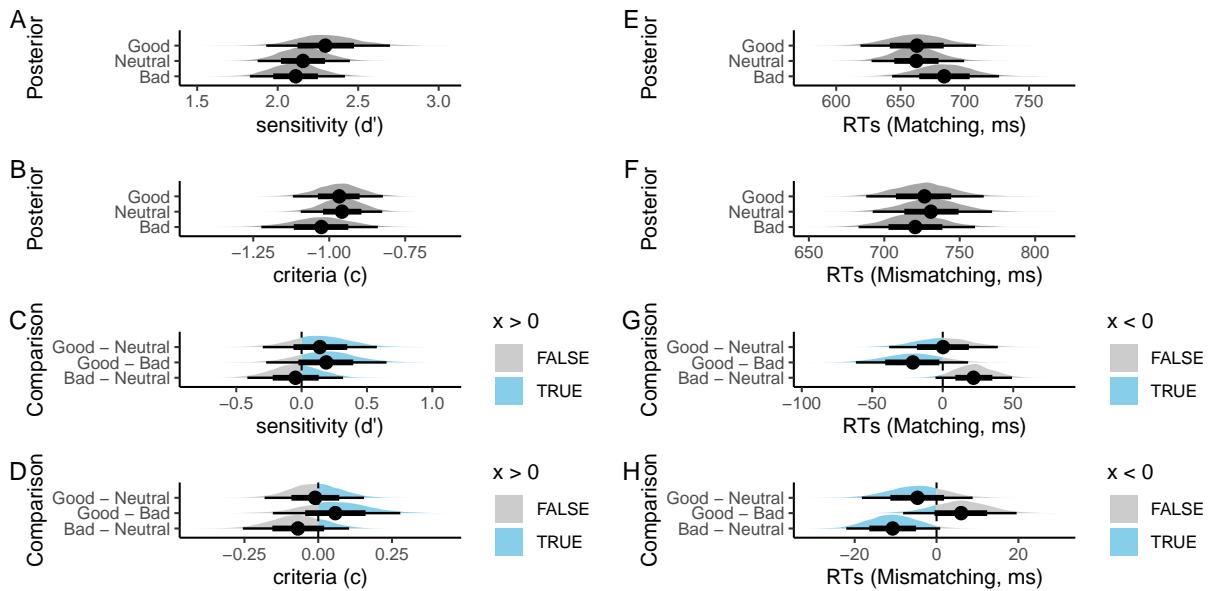


Figure 8. Exp1c: Results of Bayesian GLM analysis.

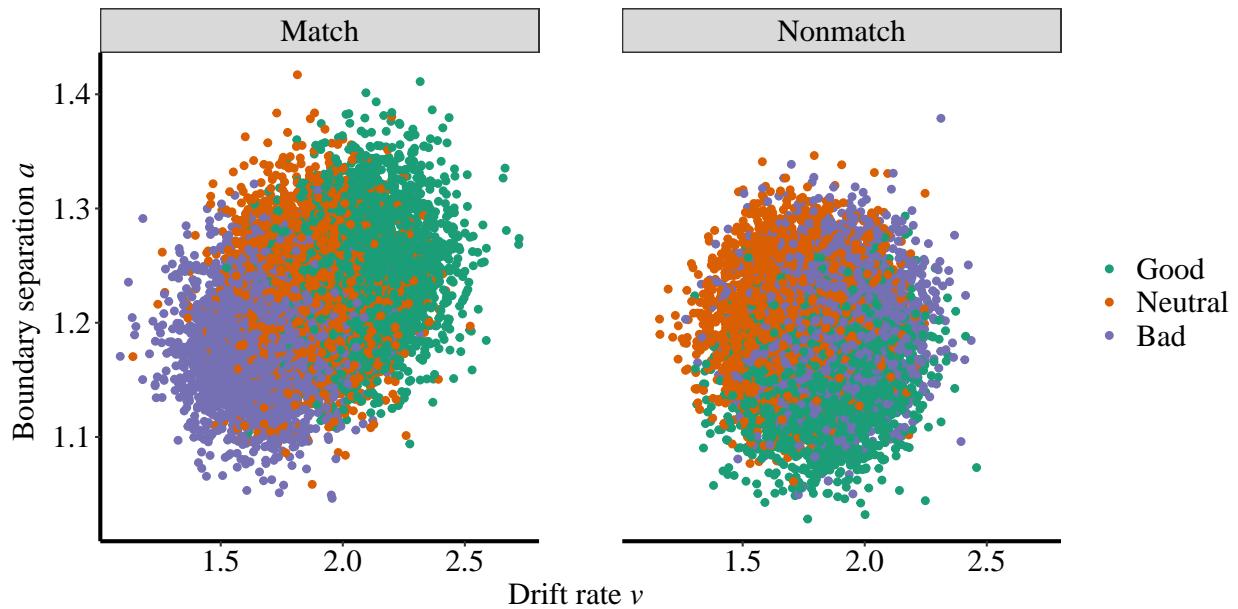


Figure 9. Exp1c: Results of HDDM.

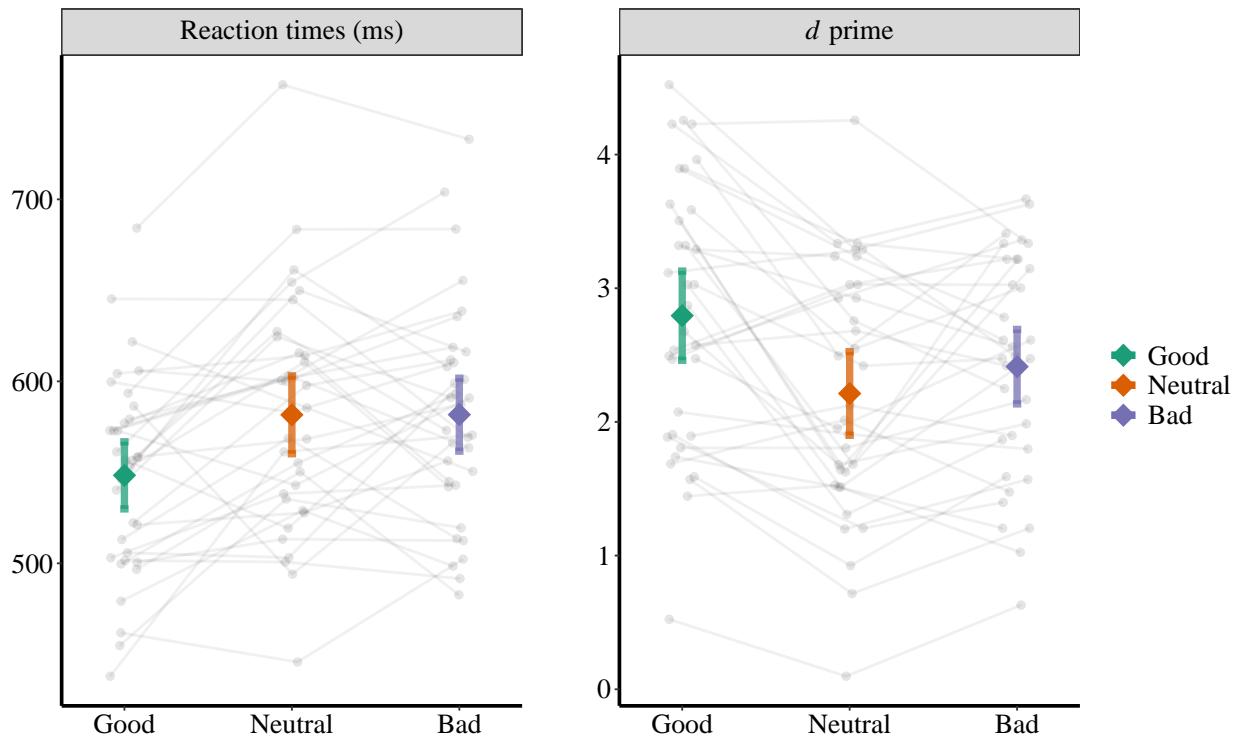
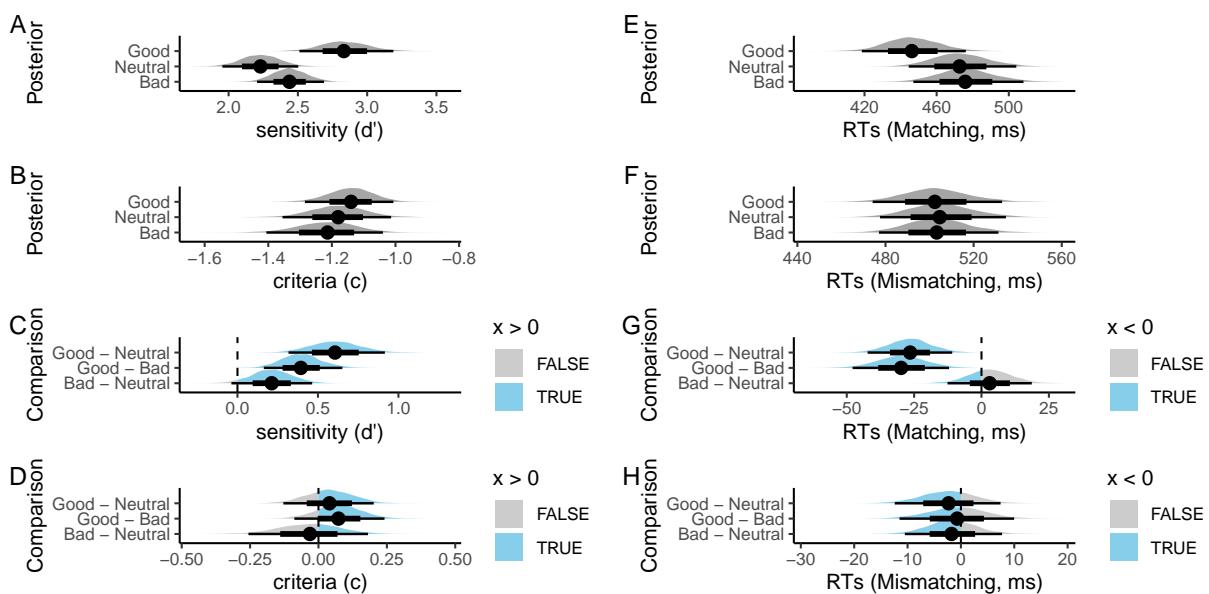
Figure 10. RT and d' prime of Experiment 2.

Figure 11. Exp2: Results of Bayesian GLM analysis.

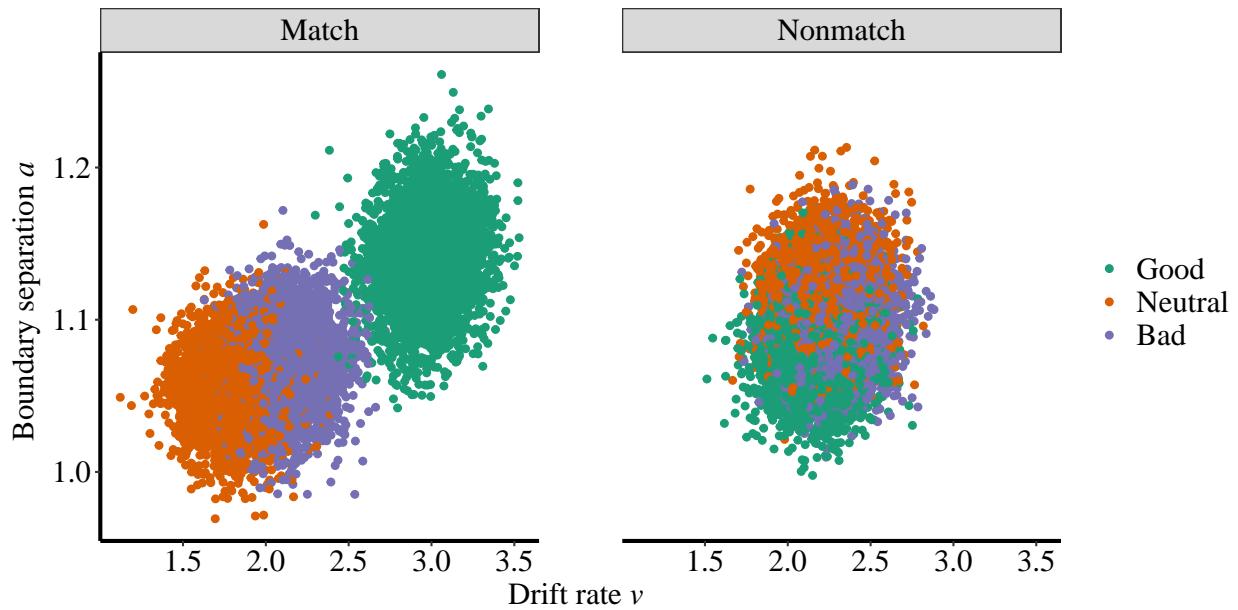


Figure 12. Exp2: Results of HDDM.

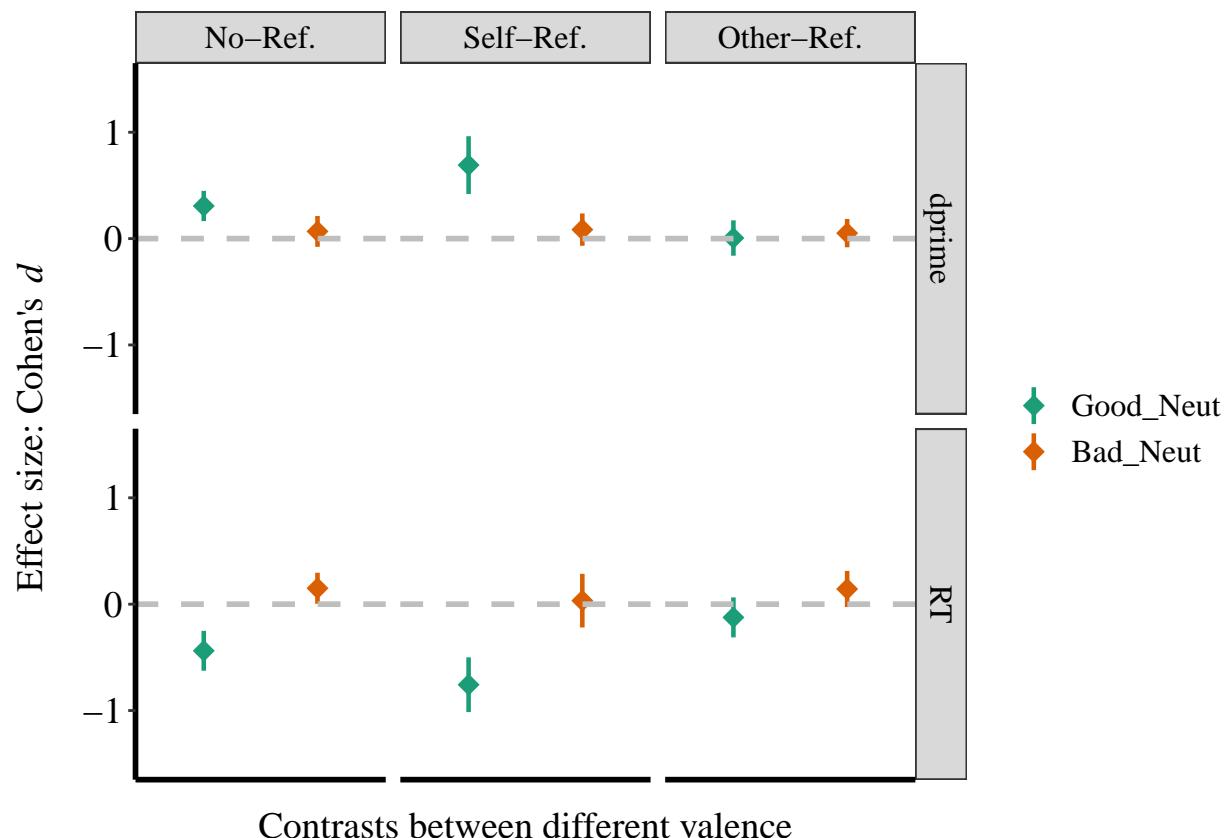


Figure 13. Effect size (Cohen's d) of Valence.

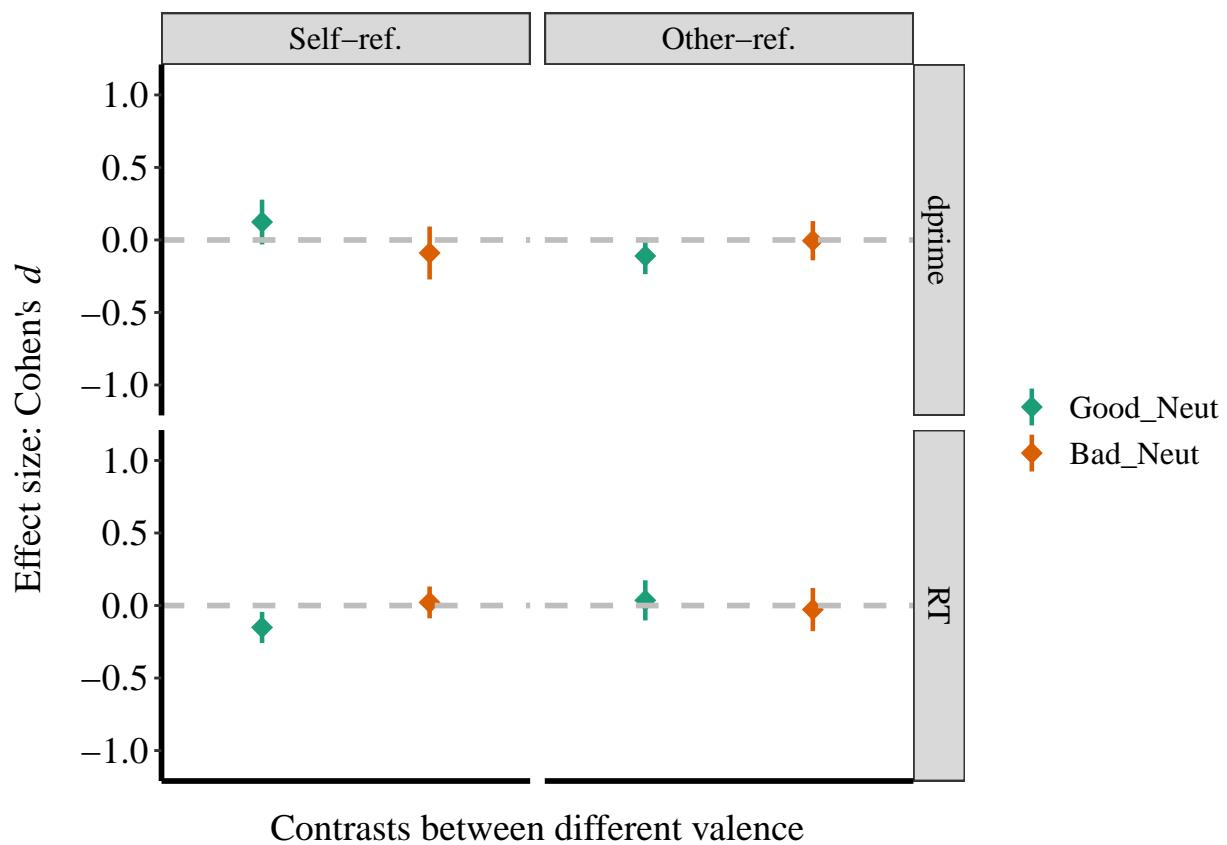


Figure 14. Effect size (Cohen's d) of Valence in Exp4a.

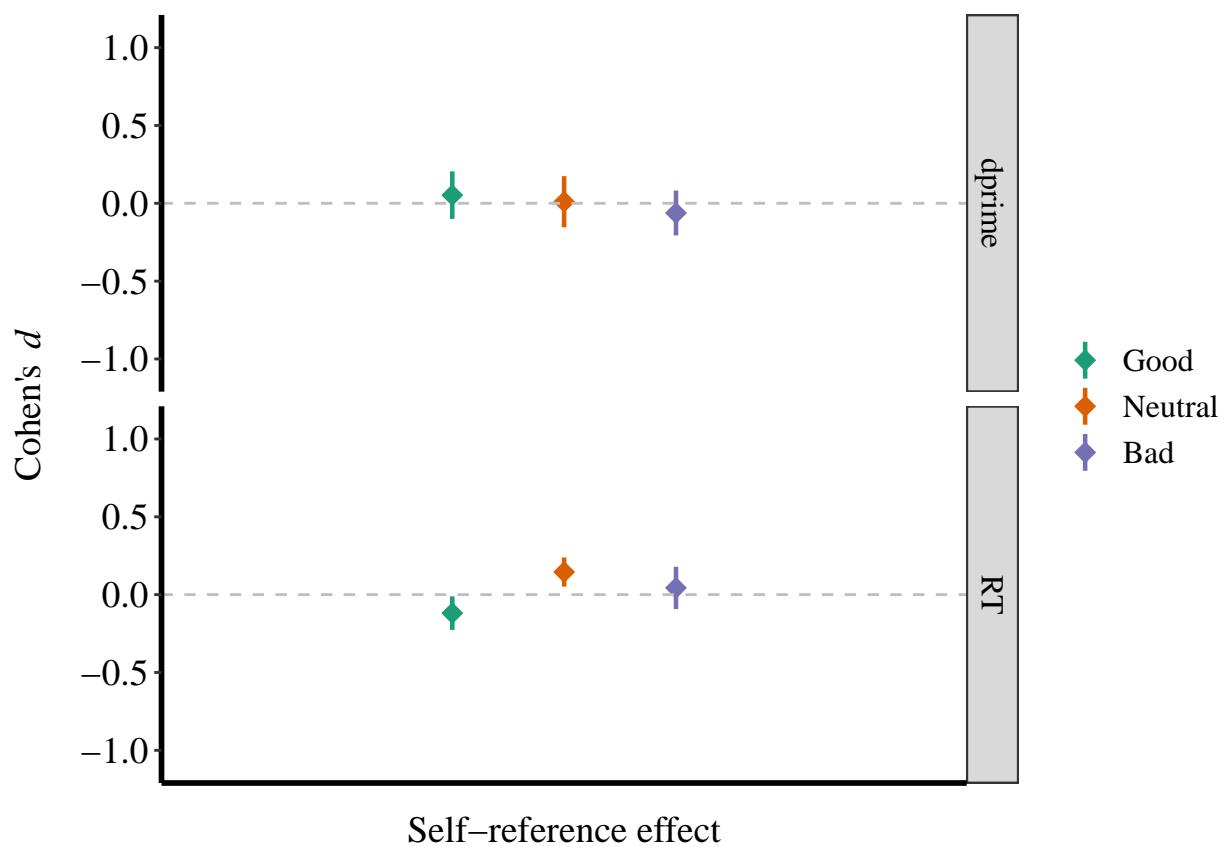


Figure 15. Effect size (Cohen's d) of Valence in Exp4b.

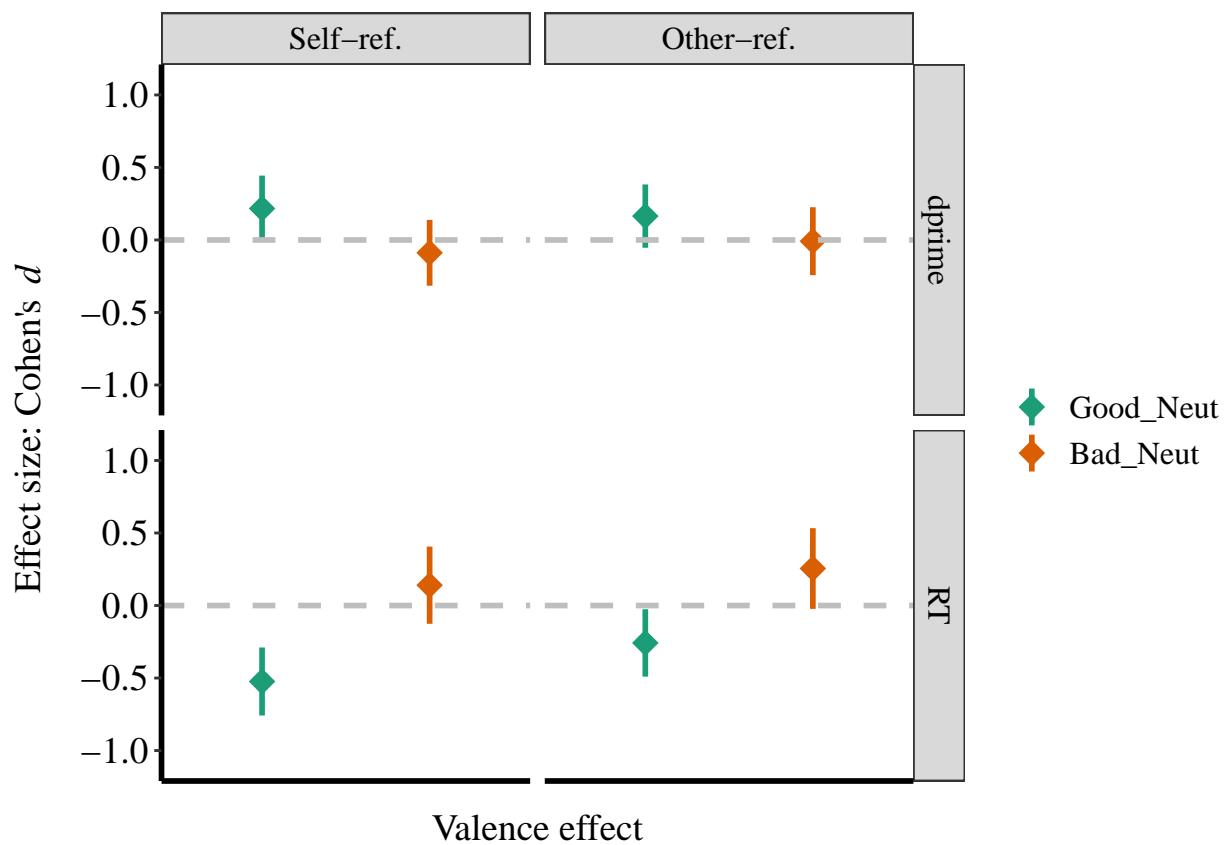


Figure 16. Effect size (Cohen's d) of Valence in Exp4b.

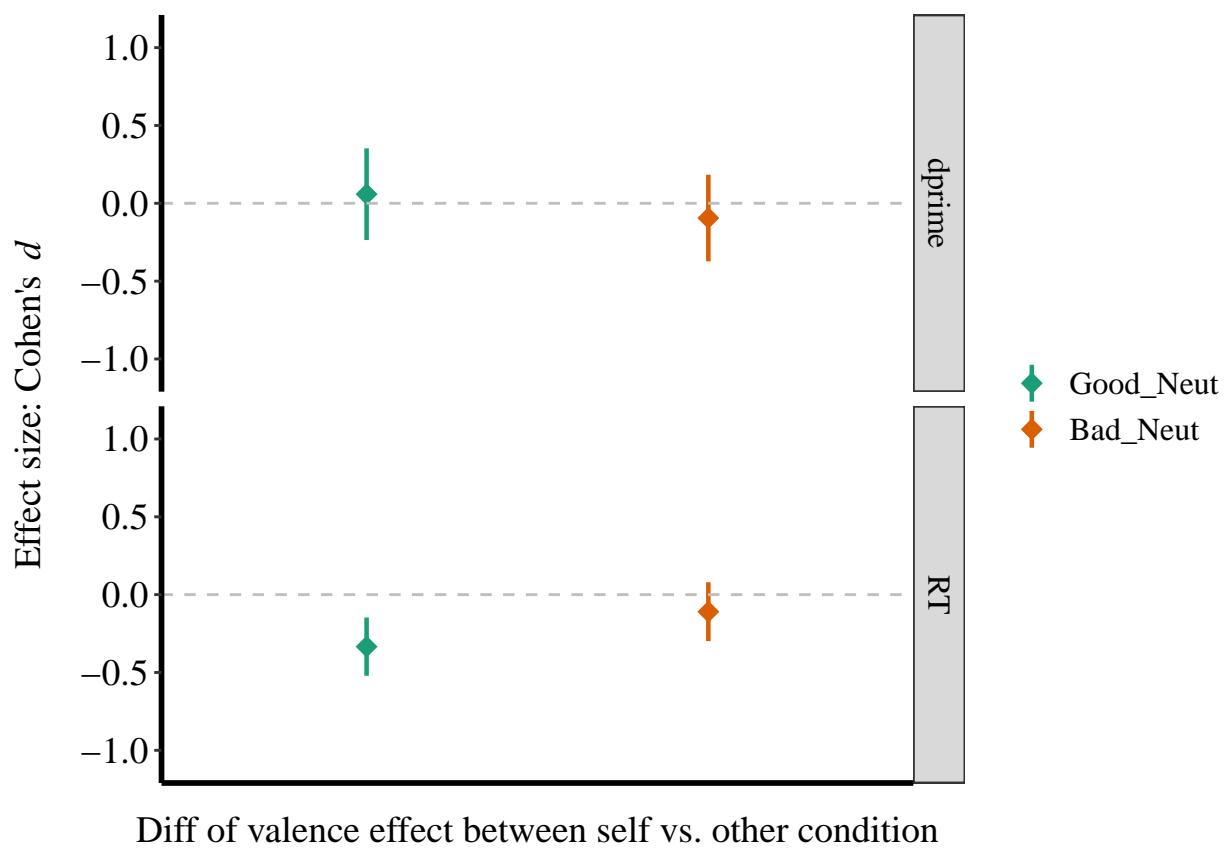


Figure 17. Effect size (Cohen's d) of Valence in Exp4b.

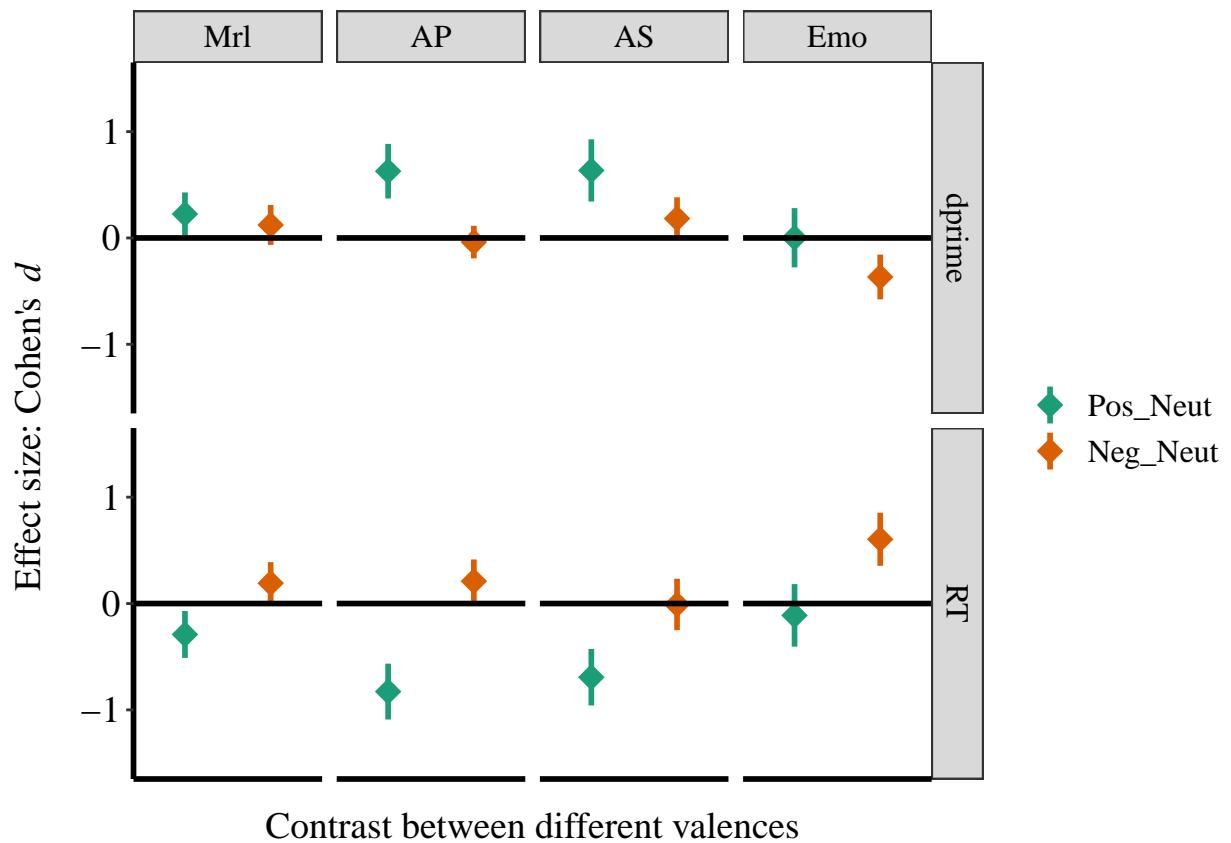


Figure 18. Effect size (Cohen's d) of Valence in Exp5.

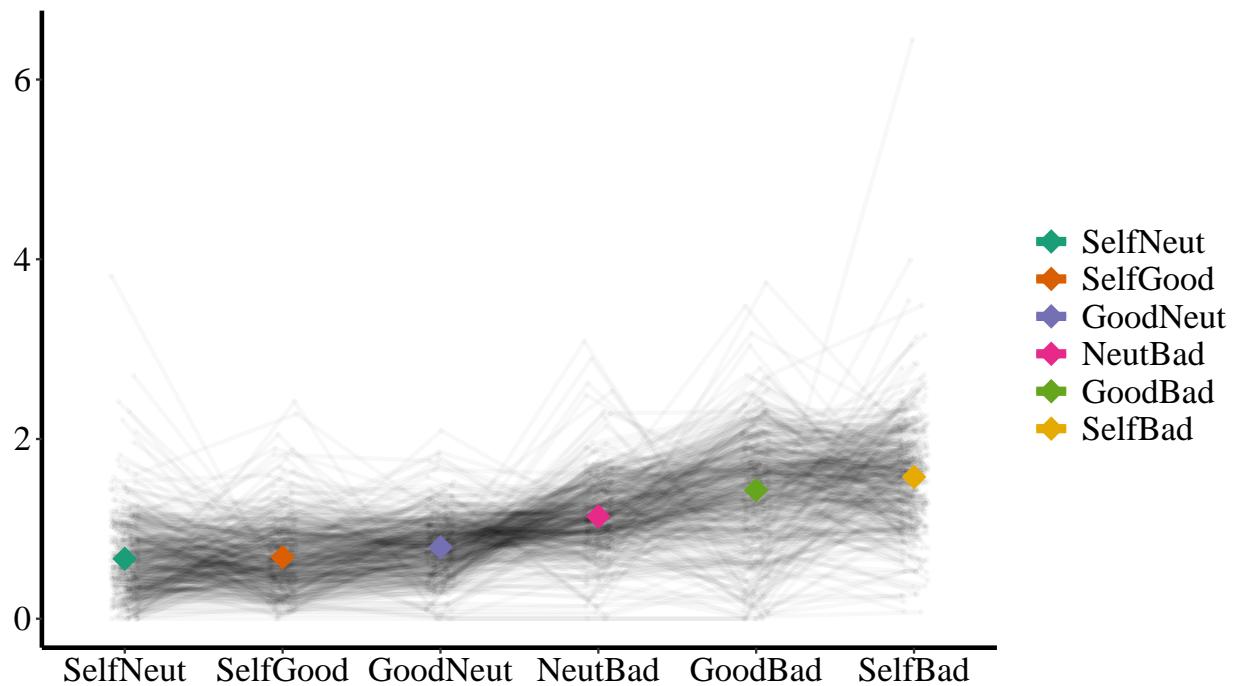


Figure 19. Self-rated personal distance

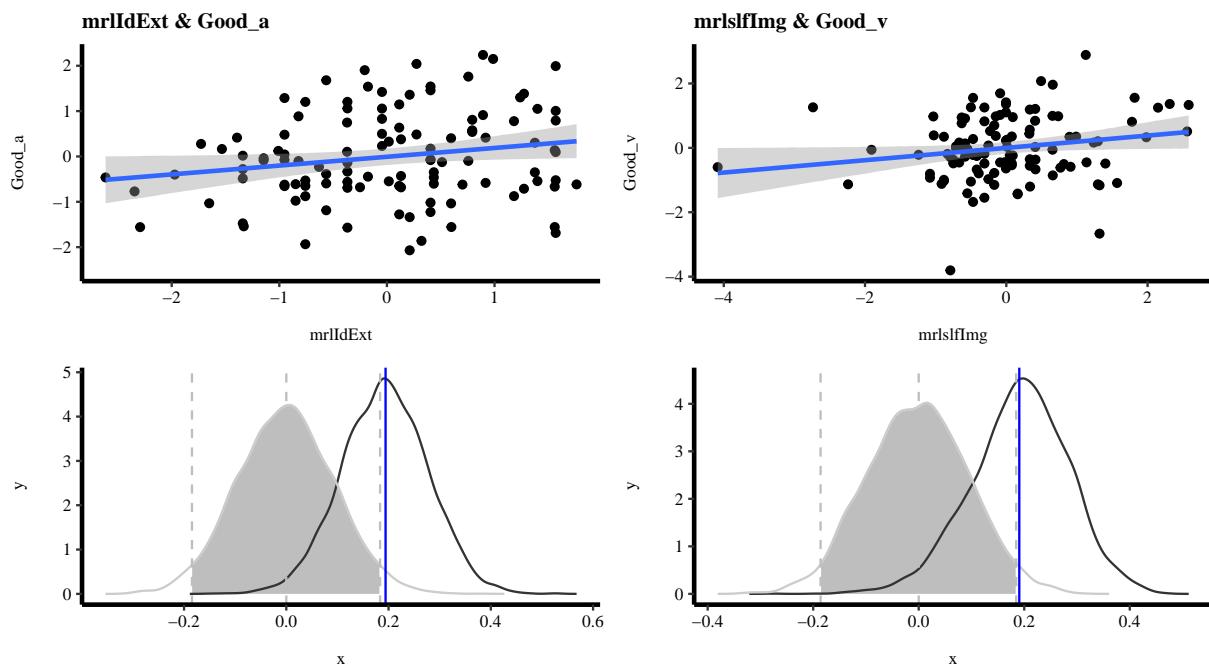


Figure 20. Correlation between moral identity and boundary separation of good condition; moral self-image and drift rate of good condition

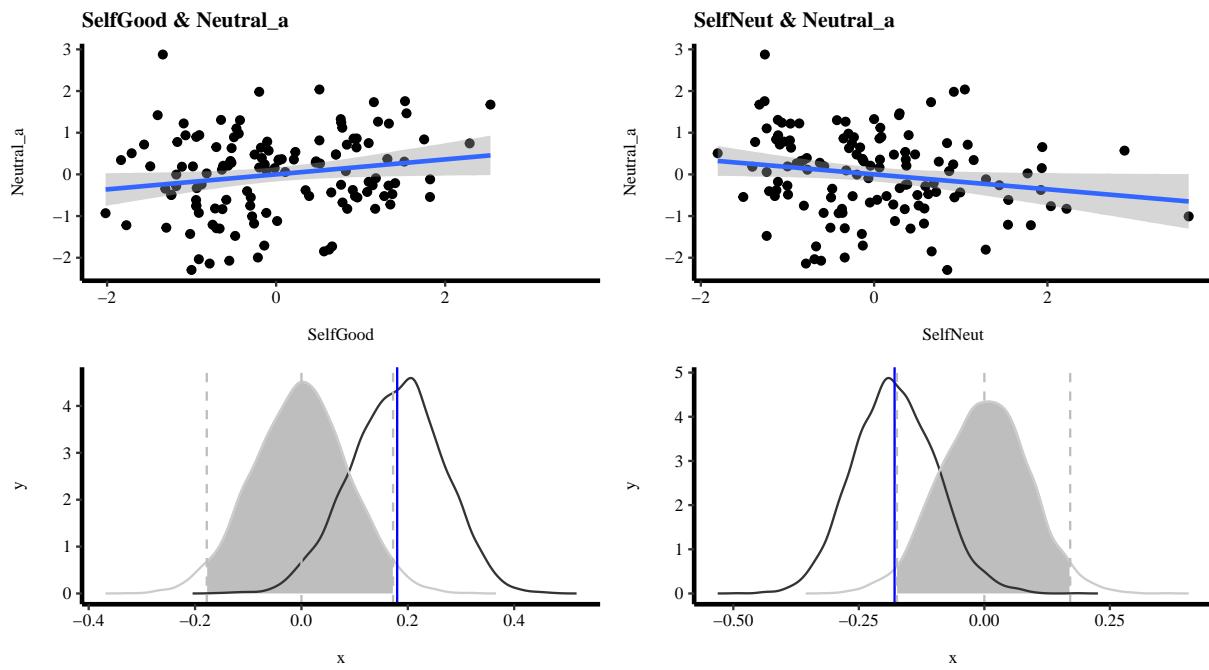


Figure 21. Correlation between personal distance and boundary separation of neutral condition

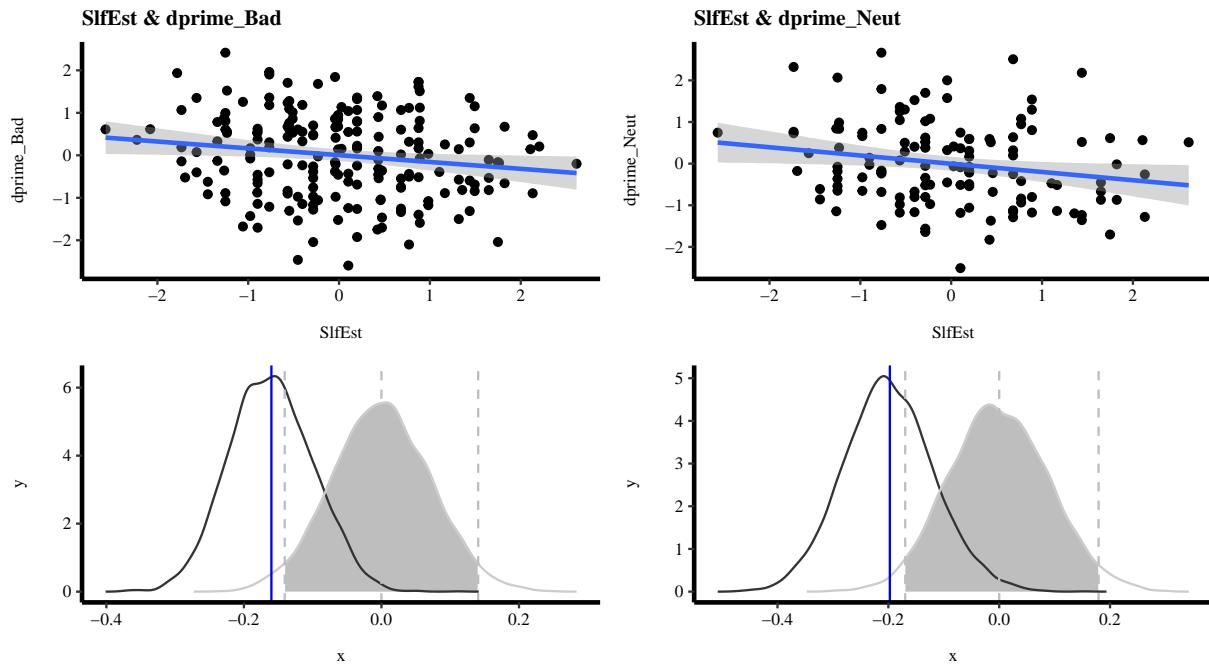


Figure 22. Correlation between self esteem and d prime of bad and neutral conditions