1                    Self-referencing prioritizes moral character on perceptual matching

Abstract

Evidence for the prioritization of moral information in cognitive processes is mixed. We examined this question using a series of eleven experiments where participants first learned associations between moral characters and geometric shapes and then performed simple speed tasks. In the first six experiments, we tested and validated prioritized responses to good characters over bad and neutral characters. To pin down the processes that are critical to the prioritization effects, in the remaining five experiments, we examined two opposing hypotheses: the valence hypothesis suggests that a general positivity bias towards all underpins the effects, while the self-binding account posits that self-referencing, rather than other-referencing is the fundamental driver of the effects. The data support the latter. Together, these results show a robust prioritization effect of good character through self-referencing processes, indicating the innate connection between morality and oneself and how humans use self-reference to explore the world and learn morality.

*Keywords:* Perceptual matching, self positivity bias, primacy of morality, Bayesian hierarchical models

Word count: X

<sub>18</sub>              Self-referencing prioritizes moral character on perceptual matching

<sub>19</sub>                                        **Introduction**

<sub>20</sub>        Morality is central to human life (Haidt & Kesebir, 2010). Thus, gathering

<sub>21</sub>  information about morality efficiently and accurately is crucial for individuals to navigate

<sub>22</sub>  the social world (Brambilla, Sacchi, Rusconi, & Goodwin, 2021). The importance of

<sub>23</sub>  morality naturally leads to the hypothesis that morality-related information is prioritized in

<sub>24</sub>  information processing, especially when attentional resources are limited. This hypothesis

<sub>25</sub>  is plausible because a large volume of studies has reported that valuable stimuli are

<sub>26</sub>  prioritized, e.g., threatening stimuli (e.g., Ohman, Lundqvist, & Esteves, 2001), rewards

<sub>27</sub>  (B. A. Anderson, Laurent, & Yantis, 2011; Sui & Humphreys, 2015a), or self-related stimuli

<sub>28</sub>  (Sui & Rotshtein, 2019). Consistent with this hypothesis, a few studies reported a

<sub>29</sub>  prioritization effect of negative moral information in visual processing: negative moral trait

<sub>30</sub>  words (Fiske, 1980; Gantman & Van Bavel, 2014; Ybarra, Chan, & Park, 2001) and faces

<sub>31</sub>  associated with bad behaviors (E. Anderson, Siegel, Bliss-Moreau, & Barrett, 2011;

<sub>32</sub>  Eiserbeck & Abdel Rahman, 2020) attracted more attention and were responded faster.

<sub>33</sub>        However, evidence for this negative moral bias effect is mixed. First, the opposite

<sub>34</sub>  effect was also reported. For example, Shore and Heerey (2013) found that faces with

<sub>35</sub>  positive interaction in a trust game were prioritized in the pre-attentive process. Also,

<sub>36</sub>  Abele and Bruckmueller found faster responses to moral words were not moderated by

<sub>37</sub>  valence (Abele & Bruckmüller, 2011). Second, the robustness of the negative moral bias

<sub>38</sub>  effect is questioned, a direct replication study failed to support the conclusion that faces

<sub>39</sub>  associated with bad social behaviors dominate visual awareness (eg., Stein, Grubb,

<sub>40</sub>  Bertrand, Suh, & Verosky, 2017). Third, the prioritization effect of morality might be

<sub>41</sub>  confounded with other factors, such as the priming effect (Firestone & Scholl, 2015, 2016b;

<sub>42</sub>  Jussim, Crawford, Anglin, Stevens, & Duarte, 2016) or differences between lexical

<sub>43</sub>  characteristics (Larsen, Mercer, & Balota, 2006). As a result, while the importance of

44 morality is widely recognized and there is initial evidence for a negative moral bias,

45 whether moral information is prioritized in perceptual processing is still an open question.

46      Here, we conducted a series of well-controlled experiments to examine the

47 prioritization effect of morality and its potential mechanisms. To eliminate the priming

48 effect and other potential confounding factors, we employed a task where participants first

49 acquired moral meanings of geometric shapes during the instruction phase and then

50 performed a simple perceptual matching task. The instruction-based associative learning

51 task is based on the fact that humans can rapidly learn based on verbal instructions (e.g.,

52 Cole, Braver, & Meiran, 2017). This instruction-based associative learning task is widely

53 used in aversive learning, value-based learning, and other tasks (Atlas, 2023; Cole et al.,

54 2017; Deltomme, Mertens, Tibboel, & Braem, 2018). Unlike previous studies relies on faces

55 or words (e.g., Bortolon & Raffard, 2018; Yaoi, Osaka, & Osaka, 2021), stimuli in the

56 current study are geometric shapes, whose moral meanings were acquired right before the

57 perceptual matching task. By counter-balancing associations between shapes and valence

58 of moral characters across different participants, we controlled the effect of these shapes on

59 the matching task. Also, in the matching task, we repeatedly present a few pairs of shapes

60 and labels to participants, the results can not be explained by semantic priming

61 (Unkelbach, Alves, & Koch, 2020), which is the center of the debate on previous results

62 (Firestone & Scholl, 2015, 2016a; Gantman & Bavel, 2015, 2016; Jussim et al., 2016).

63 Finally, we conducted a series of control experiments and established that moral content,

64 rather than other factors such as familiarity of stimuli, drove the prioritization effects.

65      To pin down the factors that are central to the prioritization effects, two competing

66 hypotheses were examined. One is the valence-based account, suggesting that a general

67 positivity bias towards all underpins the prioritization effects. In fact, the account has been

68 applied to explain not only positivity biases but also negativity biases. For example, the

69 negative bias toward moral information was explained by a threat detection mechanism

70 which might be general for all negative information (e.g., Fiske, 1980). The positive bias

toward moral information, on the other hand, was explained by the positive valence of the stimuli because the stimuli imply potential benefits (Shore & Heerey, 2013). However, these explanations often ignore the fact that valence is subjective *per se* (Juechems & Summerfield, 2019). That is, being related to a person is the premise of a stimulus or outcome being of value to the person. The subjective value is "a broader concept that refers to the personal significance or importance that a person assigns to a particular stimulus or outcome" and when the outcome is affective or emotional, researchers refer to it as "valence", i.e., positive or negative (Carruthers, 2021). The subjectivity of valence leads to an alternative explanation: self-binding account (Sui & Humphreys, 2015b). The self-binding account suggests that merely associating with the self can prioritize stimuli in perception, attention, working memory, and long-term memory (Sui & Humphreys, 2015b; Sui & Rotshtein, 2019), especially for positive information (Hu, Lan, Macrae, & Sui, 2020). According to the self-binding account, the prioritization of good character is a result of spontaneous self-referencing.

        To test the valence account and self-binding account in the prioritization effect of good character, we manipulated self-relevance and instructed participants on which moral character is self-referencing and which is not. We then tested whether the prioritization of moral character is by valence or by the associations between self-relevance and moral valence. The results revealed that the prioritization effect only occurred when shapes of good characters referred to the self of participants. We confirmed these results in the subsequent experiments, where shapes of good characters did not explicitly refer to the self or others but were merely presented with labels of the self or others. Together, these data revealed a mutual facilitation effect of good character and the self, suggesting a spontaneous self-referential process as a novel mechanism underlying the prioritization of good character in perceptual matching.

<sup>96</sup> **Disclosures**

<sup>97</sup> We reported all the measurements, analyses, and results in all the experiments in the

<sup>98</sup> current study. Participants whose overall accuracy was lower than 60% were excluded from

<sup>99</sup> analyses. Also, accurate responses with less than 200ms reaction times were excluded from

<sup>100</sup> the analysis. These excluded data can be found in the shared raw data files (see

<sup>101</sup> https://doi.org/10.5281/zenodo.8031086).

<sup>102</sup> All the experiments reported were not pre-registered. Most experiments (1a ~ 4b,

<sup>103</sup> except experiment 3b) reported in the current study were first finished between 2013 to

<sup>104</sup> 2016 at Tsinghua University, Beijing, China. Participants in these experiments were

<sup>105</sup> recruited from the local community. To increase the sample size of experiments to 50 or

<sup>106</sup> more (Simmons, Nelson, & Simonsohn, 2013), we recruited additional participants from

<sup>107</sup> Wenzhou University, Wenzhou, China, in 2017 for experiments 1a, 1b, 4a, and 4b.

<sup>108</sup> Experiment 3b was finished at Wenzhou University in 2017 (See Table 1 for an overview of

<sup>109</sup> these experiments).

<sup>110</sup> All participants received informed consent and were compensated for their time.

<sup>111</sup> These experiments were approved by the ethics board in the Department of Psychology,

<sup>112</sup> Tsinghua University.

<sup>113</sup> **General methods**

<sup>114</sup> **Design and Procedure**

<sup>115</sup> This series of experiments used the perceptual matching paradigm (or self-tagging

<sup>116</sup> paradigm, see Sui, He, and Humphreys (2012)), in which participants first learned the

<sup>117</sup> associations between geometric shapes and labels of different moral characters (e.g., in the

<sup>118</sup> first three studies, the triangle, square, and circle for shapes and Chinese words for "good

<sup>119</sup> person", "neutral person", and "bad person", respectively). The associations of shapes and

₁₂₀ labels were counterbalanced across participants. The paradigm consists of a brief learning

₁₂₁ stage and a test stage. During the learning stage, participants were instructed about the

₁₂₂ association between shapes and labels. Participants started the test stage with a practice

₁₂₃ phase to familiarize themselves with the task, in which they viewed one of the shapes above

₁₂₄ the fixation while one of the labels below the fixation and judged whether the shape and

₁₂₅ the label matched the association they learned. If the overall accuracy reached 60% or

₁₂₆ higher at the end of the practicing session, participants proceeded to the experimental task

₁₂₇ of the test stage. Otherwise, they finished another practices sessions until the overall

₁₂₈ accuracy was equal to or greater than 60%. The experimental task shared the same trial

₁₂₉ structure as in the practice.

₁₃₀     Experiments 1a, 1b, 1c, 2, 5, and 6a were designed to explore and confirm the effect

₁₃₁ of moral character on perceptual matching. All these experiments shared a 2 (matching:

₁₃₂ match vs. mismatch) by 3 (moral character: good vs. neutral vs. bad person)

₁₃₃ within-subject design. Experiment 1a was the first one of the whole series of studies, which

₁₃₄ aimed to examine the prioritization of moral character and found that shapes associated

₁₃₅ with good character were prioritized. Experiments 1b, 1c, and 2 were to confirm that it is

₁₃₆ the moral character that caused the effect. More specifically, experiment 1b used different

₁₃₇ Chinese words as labels to test whether the effect was contaminated by familiarity.

₁₃₈ Experiment 1c manipulated the moral character indirectly: participants first learned to

₁₃₉ associate different moral behaviors with different Chinese names, after remembering the

₁₄₀ association, they then associated the names with different shapes and finished the

₁₄₁ perceptual matching task. Experiment 2 further tested whether the way we presented the

₁₄₂ stimuli influenced the prioritization of moral character, by sequentially presenting labels

₁₄₃ and shapes instead of simultaneous presentation. Note that a few participants in

₁₄₄ Experiment 2 also participated in Experiment 1a because we originally planned a

₁₄₅ cross-task comparison. Experiment 5 was designed to compare the prioritization of good

₁₄₆ character with other important social values (aesthetics and emotion). All social values

had three levels, positive, neutral, and negative, and were associated with different shapes. Participants finished the associative learning task for different social values in different blocks, and the order of the social values was counterbalanced. Only the data from moral character blocks, which shared the design of experiment 1a, were reported here. Experiment 6a, which shared the same design as Experiment 2, was an EEG experiment aimed at exploring the neural mechanism of the prioritization of good character. Only behavioral results of Experiment 6a were reported here.

Experiments 3a, 3b, and 6b were designed to test whether the prioritization of good character can be explained by the valence account or by the self-binding account. For this purpose, we included self-reference as another within-subject variable. For example, Experiment 3a extended Experiment 1a into a 2 (matching: match vs.mismatch) by 2 (reference: self vs. other) by 3 (moral character: good vs. neutral vs. bad) within-subject design. Thus, in Experiment 3a, there were six conditions (good-self, neutral-self, bad-self, good-other, neutral-other, and bad-other) and six shapes (triangle, square, circle, diamond, pentagon, and trapezoids). Experiment 6b was an EEG experiment based on Experiment 3a but presented the label and shape sequentially. Because of the relatively high working memory load (six label-shape pairs), participants finished Experiment 6b in two days. On the first day, participants completed the perceptual matching task as a practice, and on the second day, they finished the task again while the EEG signals were recorded. We only focus on the first day's data here. Experiment 3b was designed to test whether the effect found in Experiments 3a and 6b is robust if we separately present the self-referencing trials and other-referencing trials. That is, participants finished two types of blocks: in the self-referencing blocks, they only made matching judgments to shape-label pairs that related to the self (i.e., shapes and labels of good-self, neutral-self, and bad-self), in the other-referencing blocks, they only responded to shape-label pairs that related to the other (i.e., shapes and labels of good-other, neutral-other, and bad-other).

Experiments 4a and 4b were designed to test whether the self and the good character

bind spontaneously. In Experiment 4a, participants were instructed to learn the association between two shapes (circle and square) with two labels (self vs. other) in the learning stage. In the test stage, they were instructed only respond to the shape and label during the test stage. However, we presented the labels of different moral characters in the shapes and instructed participants to ignore these labels when making matching judgments. If the self and good character bind together spontaneously, then the mere presence of good character will facilitate the response to shapes associated with the self. In the Experiment 4b, we reversed the role of self and moral character in the task: Participants learned associations between three moral labels (good-person, neutral-person, and bad-person) and three shapes (circle, square, and triangle) and made matching judgments about the shape and label of moral character, while words related to identity, "self" or "other", were presented within the shapes. As in Experiment 4a, participants were told to ignore the words inside the shape during the perceptual matching task. In the same vein, if the self and good character bind together spontaneously, then the mere presence of the self will facilitate the response to shapes associated with good character.

**Stimuli and Materials**

We used E-prime 2.0 for presenting stimuli and collecting behavioral responses. Data were collected from two universities located in two different cities in China. Participants recruited from Tsinghua University, Beijing, finished the experiment individually in a dim-lighted chamber. Stimuli were presented on 22-inch CRT monitors and participants rested their chins on a brace to fix the distance between their eyes and the screen around 60 cm. The visual angle of geometric shapes was about $3.7° \times 3.7°$, the fixation cross is of $0.8° \times 0.8°$ visual angle at the center of the screen. The words were of $3.6° \times 1.6°$ visual angle. The distance between the center of shapes or images of labels and the fixation cross was of $3.5°$ visual angle. Participants from Wenzhou University, Wenzhou, finished the experiment in a group consisting of $3 \sim 12$ participants in a dim-lighted testing room. They

200 were instructed to complete the whole experiment independently. Also, they were told to

201 start the experiment at the same time so that the distraction between participants was

202 minimized. The stimuli were presented on 19-inch CRT monitors with the same set of

203 parameters in E-prime 2.0 as in Tsinghua University, however, the visual angles could not

204 be controlled because participants' chins were not fixed.

205     In most of these experiments, participants were also asked to fill out questionnaires

206 following the behavioral tasks. All the questionnaire data were open (see, dataset 4 in Liu

207 et al., 2020). See Table 1 for a summary of information about all the experiments.

**Data analysis**

209     We used the `tidyverse` of r (see script `Load_save_data.r`) to preprocess the data.

210 The data from all experiments were then analyzed using Bayesian hierarchical models.

211     We used the Bayesian hierarchical model (BHM, or Bayesian generalized linear mixed

212 models, Bayesian multilevel models) to model the reaction time and accuracy data because

213 BHM provided three advantages over the classic NHST approach (repeated measure

214 ANOVA or *t*-tests). First, BHM estimates the posterior distributions of parameters for

215 statistical inference, therefore providing uncertainty in estimation (Rouder & Lu, 2005).

216 Second, BHM, where generalized linear mixed models could be easily implemented, can use

217 distributions that fit the data, instead of using the normal distribution for all data. Using

218 appropriate distributions for the data will avoid misleading results and provide a better

219 fitting of the data. For example, Reaction times are not normally distributed but are often

220 right skewed, and the linear assumption in ANOVAs is not satisfied (Rousselet & Wilcox,

221 2020). Third, BHM provides a unified framework to analyze data from different levels and

222 different sources, avoiding information loss when we need to combine data from different

223 experiments.

224     We used the `r` package `BRMs` (Bürkner, 2017), which used Stan (Carpenter et al.,

Table 1

*Information about all experiments.*

| ExpID | Time | Location | N | n.of.trials | Self.ref | Stim.for.Morality | Presenting.order |
|-------|------|----------|---|-------------|----------|-------------------|------------------|
| Exp_1a_1 | 2014-04 | Beijing | 38 (35) | 60 | NA | words | Simultaneously |
| Exp_1a_2 | 2017-04 | Wenzhou | 18 (16) | 120 | NA | words | Simultaneously |
| Exp_1b_1 | 2014-10 | Beijing | 39 (27) | 60 | NA | words | Simultaneously |
| Exp_1b_2 | 2017-04 | Wenzhou | 33 (25) | 120 | NA | words | Simultaneously |
| Exp_1c | 2014-10 | Beijing | 23 (23) | 60 | NA | descriptions | Simultaneously |
| Exp_2 | 2014-05 | Beijing | 35 (34) | 60 | NA | words | Sequentially |
| Exp_3a | 2014-11 | Beijing | 38 (35) | 60 | explicit | words | Simultaneously |
| Exp_3b | 2017-04 | Wenzhou | 61 (56) | 60 | explicit | words | Simultaneously |
| Exp_4a_1 | 2015-06 | Beijing | 32 (29) | 30 | implicit | words | Simultaneously |
| Exp_4a_2 | 2017-04 | Wenzhou | 32 (30) | 60 | implicit | words | Simultaneously |
| Exp_4b_1 | 2015-10 | Beijing | 34 (32) | 60 | implicit | words | Simultaneously |
| Exp_4b_2 | 2017-04 | Wenzhou | 19 (13) | 60 | implicit | words | Simultaneously |
| Exp_5 | 2016-01 | Beijing | 43 (38) | 60 | NA | words | Simultaneously |
| Exp_6a | 2014-12 | Beijing | 24 (24) | 180 | NA | words | Sequentially |
| Exp_6b | 2016-01 | Beijing | 23 (22) | 90 | explicit | words | Sequentially |

*Note.* Stim.for.Morality = How moral character was manipulated; Presenting.order = How shapes & labels were presented. Number in () for N is number of participants are included in the analysis. In the current analysis, we only remain participants' data when they participate the experiment for the first time.

2017) as the back-end, for the BHM analyses. We estimated the overall effect across

experiments that shared the same experimental design using one model, instead of a

two-step approach that was adopted in mini-meta-analysis (e.g., Goh, Hall, & Rosenthal,

2016). More specifically, a three-level model was used to estimate the overall effect of

prioritization of good character, which included data from five experiments: 1a, 1b, 1c, 2,

5, and 6a. Similarly, a three-level HBM model is used for experiments 3a, 3b, and 6b.

Results of individual experiments can be found in the supplementary results. For

experiments 4a and 4b, which tested the implicit interaction between the self and good

character, we used HBM for each experiment separately.

For questionnaire data, we only reported the subjective distance between different

persons or moral characters in the supplementary results and did not analyze other

questionnaire data in the present study, which were described in (Liu et al., 2020).

**Response data.**   We followed previous studies (Hu et al., 2020; Sui et al., 2012)

and used the signal detection theory approach to analyze the response accuracy. More

specifically, the match trials are treated as signals and non-match trials are noise. The

sensitivity and criterion of signal detection theory are modeled through BHM (Rouder &

Lu, 2005).

We used the Bernoulli distribution for the signal detection theory. The probability

that the $j$th subject responded "match" ($y_{ij} = 1$) at the $i$th trial $p_{ij}$ is distributed as a

Bernoulli distribution with parameter $p_{ij}$:

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

The reparameterized value of $p_{ij}$ is a linear regression of the independent variables:

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

where the probits (z-scores; $\Phi$, "Phi") of $p$s is used for the regression.

$_{247}$      The participant-specific intercepts ($\beta_0 = -zFAR$) and slopes ($\beta_1 = d'$) are described

$_{248}$  by multivariate normal with means and a covariance matrix for the parameters.

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

$_{249}$      We used the following formula for Experiments 1a, 1b, 1c, 2, 5, and 6a, which have a

$_{250}$  2 (matching: match vs. mismatch) by 3 (moral character: good vs. neutral vs. bad)

$_{251}$  within-subject design:

$_{252}$      `saymatch ~ 0 + Valence + Valence:ismatch + (0 + Valence +`

$_{253}$  `Valence:ismatch | Subject) + (0 + Valence + Valence:ismatch |`

$_{254}$  `ExpID_new:Subject), family = bernoulli(link="probit")`

$_{255}$      in which the `saymatch` is the response data whether participants pressed the key

$_{256}$  corresponding to "match", `mismatch` is the independent variable of matching, `Valence` is

$_{257}$  the independent variable of moral character, `Subject` is the index of participants, and

$_{258}$  `Exp_ID_new` is the index of different experiments. Note that we distinguished data

$_{259}$  collected from two universities.

$_{260}$      For experiments 3a, 3b, and 6b, an additional variable, i.e., reference (self vs. other),

$_{261}$  was included in the formula:

$_{262}$      `saymatch ~ 0 + ID:Valence + ID:Valence:ismatch + (0 + ID:Valence +`

$_{263}$  `ID:Valence:ismatch | Subject) + (0 + ID:Valence + ID:Valence:ismatch |`

$_{264}$  `ExpID_new:Subject), family = bernoulli(link="probit")`

$_{265}$      in which the `ID` is the independent variable "reference", which means whether the

$_{266}$  stimulus was self-referencing or other-referencing.

$_{267}$      **Reaction times.**   We used log-normal distribution to model the RT data (see

$_{268}$  https://lindeloev.github.io/shiny-rt/#34_(shifted)_log-normal). This means we need to

$_{269}$  estimate the posterior of two parameters: $\mu$, and $\sigma$. $\mu$ is the mean of the `logNormal`

$_{270}$  distribution, and $\sigma$ is the disperse of the distribution.

271   The reaction time of the $j$th subject on $i$th trial, $y_{ij}$, is log-normal distributed:

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

272   The parameter $\mu_j$ is a linear regression of the independent variables:

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

273   and the parameter $\sigma_j$ does not vary with independent variables:

$$\sigma_j \sim HalfNormal()$$

274   The participant-specific intercepts $(\beta_{0j})$ and slopes $(\beta_{1j})$ are described by
275   multivariate normal with means and a covariance matrix for the parameters.

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

276   The formula used for experiments 1a, 1b, 1c, 2, 5, and 6a, which have a 2 (matching:
277   match vs. non-match) by 3 (moral character: good vs. neutral vs. bad) within-subject
278   design, is as follows:

279   `RT_sec ~ 1 + Valence*ismatch + (Valence*ismatch | Subject) +`
280   `(Valence*ismatch | ExpID_new:Subject), family = lognormal()`

281   in which `RT_sec` is the reaction times data with the second as a unit. The other
282   variables in this formula have the same meaning as the response data.

283   For experiments 3a, 3b, and 6b, which have a 2 by 2 by 3 within-subject design, the
284   formula is as follows: `RT_sec ~ 1 + ID*Valence + (ID*Valence | Subject) +`
285   `(ID*Valence | ExpID_new:Subject), family = lognormal()` .

286   Note that for experiments 3a, 3b, and 6b, the three-level model for reaction times
287   only included the matched trials to avoid divergence when estimating the posterior of the
288   parameters.

289     **Testing hypotheses.**    To test hypotheses, we used the Sequential Effect eXistence

290 and sIgnificance Testing (SEXIT) framework suggested by Makowski, Ben-Shachar, Chen,

291 and Lüdecke (2019). In this approach, we used the posterior distributions of model

292 parameters or other effects that can be derived from posterior distributions. The SEXIT

293 approach reports centrality, uncertainty, existence, significance, and size of the input

294 posterior, which is intuitive for making statistical inferences. We used `bayestestR` for

295 implementing this approach (Makowski, Ben-Shachar, & Lüdecke, 2019).

296     ***Prioritization of moral character.***    We tested whether moral characters are

297 prioritized by examining the population-level effects (also called fixed effect) of the

298 three-level Bayesian hierarchical model of Experiments 1a, 1b, 1c, 2, 5, and 6a. More

299 specifically, we calculated the differences between the posterior distributions of the

300 good/bad character and the neutral character and then tested these posterior distributions

301 with the SEXIT approach.

302     ***Modulation of self-relevance.***    We tested the modulation effect of the

303 self-referencing process by examining the interaction between moral character and the

304 self-referencing process for the three-level Bayesian hierarchical model of Experiments 3a,

305 3b, and 6b. More specifically, we tested two possible explanations for the prioritization of

306 good character: the valence effect alone or an interaction between the valence effect and

307 self-relevance. If the former is correct, then there will be no interaction between moral

308 character and self-relevance, i.e., the prioritization effect exhibits a similar pattern for both

309 self- and other-referencing conditions. Otherwise, there will be an interaction between the

310 two factors, i.e., the prioritization effect exhibits different patterns for self- and

311 other-referencing conditions. To test the interaction, we calculated the posterior

312 distribution of the difference of difference: $(good - neutral)_{self}$ vs. $(good - neutral)_{other}$.

313 We then tested the difference of difference with SEXIT approach.

314     ***Spontaneous binding between the self and good character.***    For data from

315 Experiments 4a and 4b, we further examined whether the self-referencing process is

spontaneous (i.e., whether the good character is spontaneously bound with the self). For
Experiment 4a, if there exists a spontaneous binding between self and good character,
there should be an interaction between moral character and self-relevance. More
specifically, we tested the posterior distributions of $good_{self} - neutral_{self}$ and
$good_{other} - neutral_{other}$, as well as the difference between these differences with the
SEXIT framework. For Experiment 4b, if there exists a spontaneous binding between
self-relevance and good character, then, there will be a self-other difference for some moral
character conditions but not for other moral character conditions. More specifically, we
tested the posteriors of $good_{self} - good_{other}$, $neutral_{self} - neutral_{other}$, and
$bad_{self} - bad_{other}$ as well as the difference between them with SEXIT framework.

# Results

## Prioritization of good character

To test whether moral characters are prioritized, we modeled data from Experiments
1a, 1b, 1c, 2, 5, and 6a with three-level Bayesian hierarchical models. All these experiments
shared similar designs, with a total sample size of 192. Note that for both experiments 1a
and 1b, two datasets were collected at different time points and locations, thus we treated
them as independent samples. Here we only reported the population-level results of
three-level Bayesian models, the results of each experiment can be found in supplementary
materials.

For the $d$ prime, results from the Bayesian model revealed a robust effect of moral
character. Shapes associated with good characters ("good person", "kind person" or a
name associated with good behaviors) have higher sensitivity (median = 2.45, 95% HDI =
[2.24 2.72]) than shapes associated with neutral characters (median = 2.15, 95% HDI =
[1.92 2.45]), the difference ($median_{diff} = 0.31$, 95% HDI [0, 0.62]) has a 97.31% probability
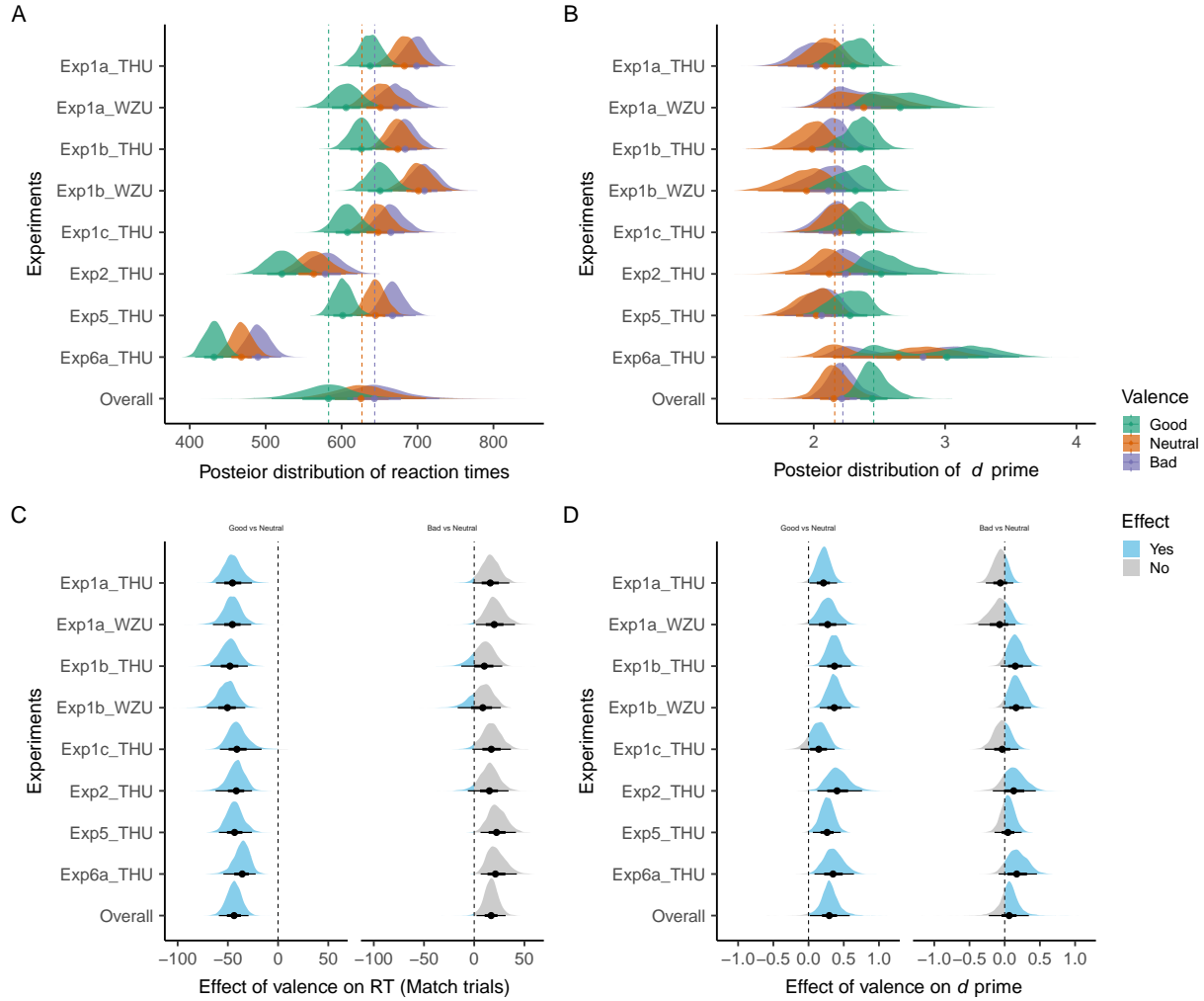of being positive ($> 0$), 94.91% of being significant ($> 0.05$). But we did not find a

341  difference between shapes associated with bad characters (median = 2.21, 95% HDI = [2.00

342  2.48]) and neutral character, the difference ($median_{diff}$ = 0.05, 95% HDI [-0.27, 0.38])

343  only has a 60.56% probability of being positive ($> 0$), 49.34% of being significant ($> 0.05$).

344        The results from reaction times also found a robust effect of moral character for both

345  match trials (see figure 1 C) and nonmatch trials (**see supplementary materials**). For

346  match trials, shapes associated with good characters were faster (median = 583 ms, 95%

347  HDI = [506 663]) than shapes associated with neutral characters (median = 626 ms, 95%

348  HDI = [547 710]), the effect ($median_{diff}$ = -44, 95% HDI [-67, -24]) has a 99.94%

349  probability of being negative ($< 0$), 99.94% of being significant ($< -0.05$). We also found

350  that RTs to shapes associated with bad characters (median = 643 ms, 95% HDI = [564

351  729]) were slower as compared to the neutral character, the effect ($median_{diff}$ = 17, 95%

352  HDI [-6, 36]) has a 93.58% probability of being positive ($> 0$), 93.55% of being significant

353  ($> 0.05$).

354        For the nonmatch trials, we found a similar pattern but a much smaller effect size.

355  Shapes associated with good characters (median = 657 ms, 95% HDI = [571 739]) were

356  faster than shapes associated with neutral characters (median = 673 ms, 95% HDI = [589

357  761]), the difference ($median_{diff}$ = -18, 95% HDI [-27, -8]) has a 99.91% probability of

358  being negative ($< 0$), 99.91% of being significant ($< -0.05$). In contrast, the shapes

359  associated with bad characters (median = 678 ms, 95% HDI = [592 764]) were slower than

360  shapes associated with neutral characters, the effect ($median_{diff}$ = 5, 95% HDI [-3, 13])

361  has a 92.43% probability of being positive ($> 0$), 92.31% of being significant ($> 0.05$).

362  **Modulation effect self-referential processing**

363        To test the modulation effect of self-relevance, we also modeled data from three

364  experiments (3a, 3b, and 6b) with three-level Bayesian models. These three experiments

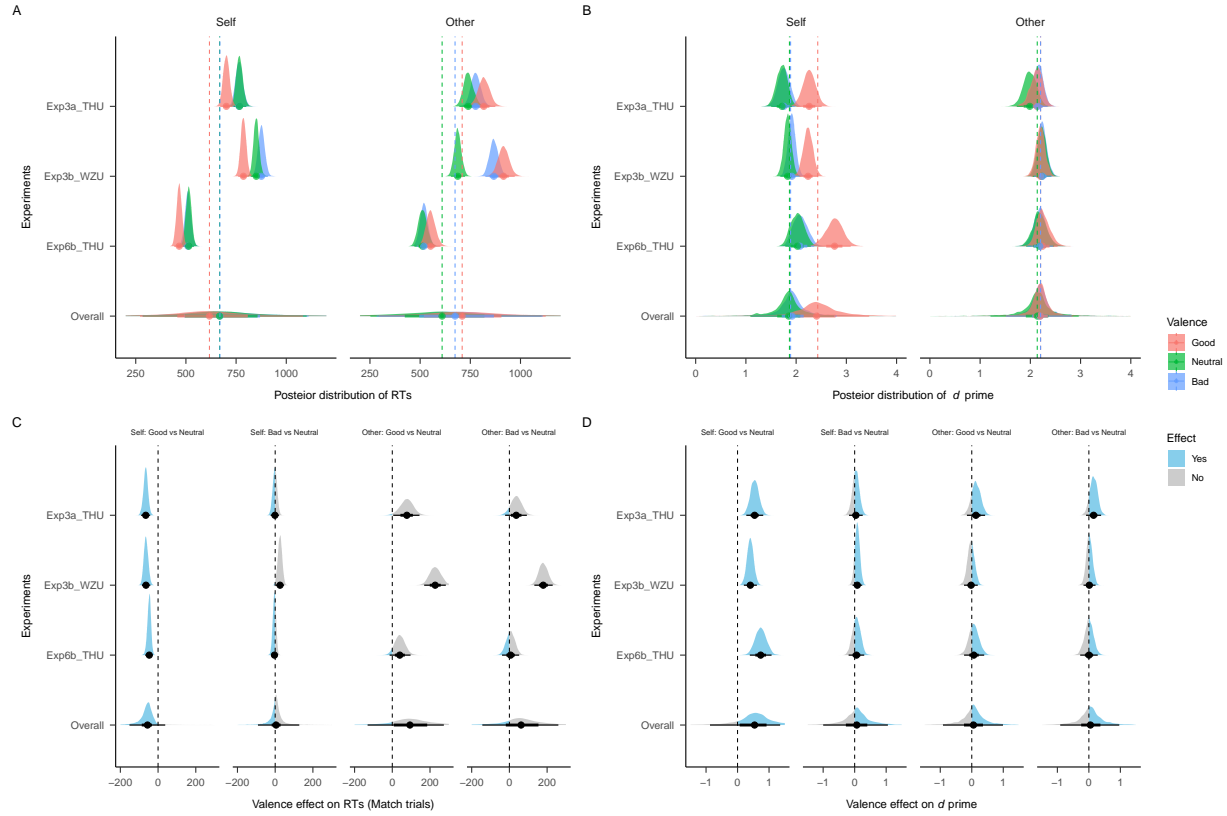365  included 108 unique participants. We focused on the population-level effect of the

*Figure 1.* Effect of moral character on perceptual matching. (A) Experimental level (six experiments, with eight independent samples) and population level posterior distributions of RT under different matching conditions; (B) Experimental level and population level posterior distributions of *d*-prime under different conditions; (C) Experimental level and population level posterior distributions of the RT differences between conditions (left, Good vs. Neutral; right, Bad vs. Neutral); (D) Experimental level and population level posterior distributions of the *d*-prime differences between conditions (left, Good vs. Neutral; right, Bad vs. Neutral).

366  interaction between self-referential processing and moral valence. Also, we examined the

367  differences of differences, i.e., how the differences between good/bad characters and the

368  neutral character under the self-referencing conditions differ from that under

369  other-referencing conditions. The results of each experiment can be found in

370  supplementary materials.

371      For the $d$ prime, we found an interaction between the moral valence and

372  self-relevance: the good-neutral differences are larger for the self-referencing condition than

373  for the other-referencing condition, the difference ($median_{diff}$ = 0.48, 95% HDI [-0.62,

374  1.65]) has a 93.04% probability of being positive ($> 0$), 91.92% of being significant ($>$

375  0.05). However, the bad-neutral differences ($median_{diff}$ = 0.0087, 95% HDI [-0.96, 1.00])

376  only have a 51.85% probability of being positive ($> 0$), 41.29% of being significant ($>$

377  0.05). Further analyses revealed that the prioritization effect of good character (as

378  compared to neutral) only appeared for self-referencing conditions but not

379  other-referencing conditions. The estimated $d$ prime for good-self was greater than

380  neutral-self ($median_{diff}$ = 0.54, 95% HDI [-0.30, 1.41]), with a 95.99% probability of being

381  positive ($> 0$), 95.36% of being significant ($> 0.05$). The differences between bad-self and

382  neutral-self, good-other and neutral-other, and bad-other and neutral-other are all centered

383  around zero (see Figure 2, B, D).

384      For the RTs of matched trials, we also found an interaction between moral valence

385  and self-relevance: the good-neutral differences were larger for the self- than the

386  other-referencing conditions ($median_{diff}$ = -148, 95% HDI [-413, 73]) has a 96.05%

387  probability of being negative ($< 0$), 96.05% of being significant ($< -0.05$). However, this

388  pattern was much weaker for bad-neutral differences ($median_{diff}$ = -47, 95% HDI [-280,

389  182]) has a 79.91% probability of being negative ($< 0$) and 79.88% of being significant ($<$

390  -0.05). Further analyses revealed a robust good-self prioritization effect as compared to

391  neutral-self ($median_{diff}$ = -59, 95% HDI [-115, -22]) has a 98.87% probability of being

392  negative ($< 0$) and 98.87% of being significant ($< -0.05$)) and good-other ($median_{diff}$ =

*Figure 2.* Interaction between moral character and self-referential. (A) Experimental level (three experiments) and population level posterior distributions of RT under different conditions; (B) Experimental level and population level posterior distributions of *d*-prime under different conditions; (C) Experimental level and population level posterior distributions of the RT differences between conditions, from left to right: Good-self vs. Neutral-self, Bad-self vs. Neutral-self, Good-other vs. Neutral-other, Bad-other vs. Neutral-other; (D) Experimental level and population level posterior distributions of the *d*-prime differences between conditions, from left to right: Good-self vs. Neutral-self, Bad-self vs. Neutral-self, Good-other vs. Neutral-other, Bad-other vs. Neutral-other.

393 -109, 95% HDI [-227, -31]) has a 98.65% probability of being negative ($< 0$) and 98.65% of

394 being significant ($< -0.05$)) conditions. Similar to the results of $d'$, we found that

395 participants responded slower for both good character than for the neutral character when

396 they referred to others, $median_{diff} = 85.01$, 95% HDI [-112, 328]) has a 92.16%

397 probability of being positive ($> 0$) and 92.15% of being significant ($> 0.05$). A similar

398 pattern was also found for the bad character when referred to others: bad-other responded

399 slower than neutral-other, $median_{diff} = 44$, 95% HDI [-146, 268]) has an 80.03%

400 probability of being positive ($> 0$) and 79.99% of being significant ($> 0.05$). See Figure 2.

401      These results suggested that the prioritization of good character is not solely driven

402 by the valence of moral character. Instead, self-relevance modulated the prioritization of

403 good character: good character was prioritized only when it referred to the self. When the

404 moral character referred to others, responses to both good and bad characters were slowed

405 down.

406 **The link between oneself and good character**

407      Experiments 4a and 4b were designed to test whether the good character and the self

408 bind together spontaneously. Because these two experiments have different experimental

409 designs, we model their data separately.

410      In experiment 4a, where "self" vs. "other" were task-relevant and moral character

411 were task-irrelevant, we found the "self" conditions performed better than the "other"

412 conditions for both $d$ prime and reaction times. This pattern is consistent with previous

413 studies (e.g., Sui et al. (2012)).

414      More importantly, we found evidence that task-irrelevant moral character also played

415 a role. For shapes associated with "self", $d'$ was greater when shapes had a good character

416 inside (median = 2.82, 95% HDI [2.64 3.03]) than shapes that have neutral character

417 (median = 2.74, 95% HDI [2.58 2.94]), the difference (median = 0.08, 95% HDI [-0.10,

0.27]) has an 81.60% probability of being positive ($> 0$), 64.33% of being significant ($>$

0.05). For shapes associated with "other", the pattern reversed: $d$ prime was smaller when

shapes had a good character inside (median = 1.87, 95% HDI [1.70 2.04]) than had neutral

(median = 1.96, 95% HDI [1.79 2.14]), the difference (median = -0.09, 95% HDI [-0.25,

0.05]) has an 89.03% probability of being negative ($< 0$), 71.38% of being significant ($<$

-0.05). The difference between these two effects (median = 0.18, 95% HDI [-0.06, 0.43]) has

a 92.88% probability of being positive ($> 0$), 85.08% being significant ($> 0.05$). See Figure

3.

A similar but more robust pattern was found for RTs in matched trials. For the "self"

condition, when a good character was presented inside the shapes, the RTs (median = 633,

95% HDI [614 654]) were faster than when a neutral character (median = 647, 95% HDI

[628 666]) was inside, the effect (median = -8, 95% HDI [-17, 2]) has a 94.55% probability

of being negative ($< 0$) and 94.50% of being significant ($< -0.05$). In contrast, when the

shapes referred to other, RTs for shapes with good character inside (median = 733, 95%

HDI [707 756]) were slower than those with neutral character inside (median = 713, 95%

HDI [691 734]), the effect (median = 12, 95% HDI [-4, 28]) has a 93.00% probability of

being positive ($> 0$) and 92.83% of being significant ($> 0.05$). The difference between these

effects (median = -19, 95% HDI [-43, 4]) has a 94.90% probability of being negative ($< 0$)

and 94.88% of being significant ($< -0.05$).

In experiment 4b, where moral characters were task-relevant and "self" vs "other"

were task-irrelevant, we found a main effect of moral character: performance for shapes

associated with good characters was better than other-related conditions on both *d'* and

reaction times. This pattern, again, shows a robust prioritization effect of good character.
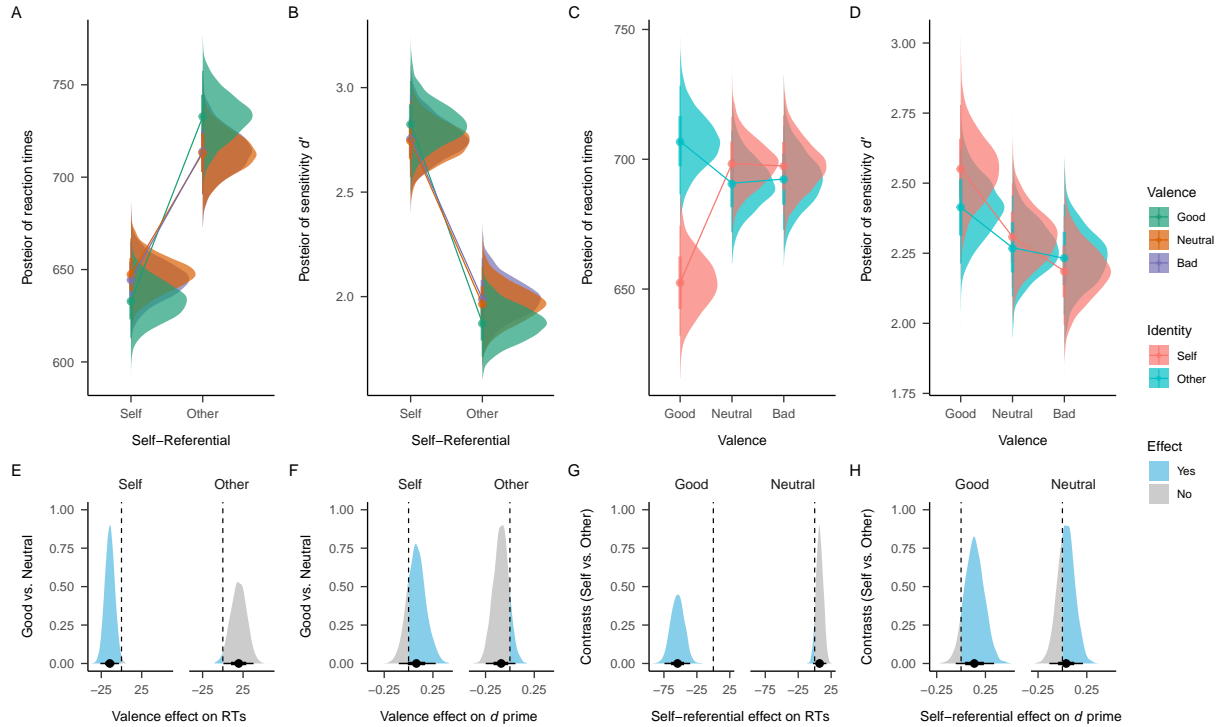
Most importantly, we found evidence that task-irrelevant labels, "self" or "other",

also played a role. For shapes associated with good character, the $d$ prime was greater

when shapes had a "self" inside than with "other" inside ($mean_{diff} = 0.14$, 95% HDI

[-0.05, 0.34]) has a 92.35% probability of being positive ($> 0$) and 81.80% of being

significant ($> 0.05$). However, the difference did not occur when the target shape where

associated with "neutral" ($mean_{diff} = 0.04$, 95% HDI [-0.13, 0.22]) and has a 67.20%

probability of being positive ($> 0$) and 44.80% of being significant ($> 0.05$). Neither for the

"bad" person condition: $mean_{diff} = 0.10$, 95% HDI [-0.16, 0.37]) has a 77.03% probability

of being positive ($> 0$) and 64.62% of being significant ($> 0.05$).

The same trend appeared for the RT data. For shapes associated with good

character, having a "self" inside shapes reduced the reaction times as compared to having

an "other" inside the shapes ($mean_{diff} = -55$, 95% HDI [-75, -35]) has a 100% probability

of being negative ($< 0$) and 100.00% of being significant ($< -0.05$). However, when the

shapes were associated with the neutral character, having a "self" inside shapes increased

the RTs: $mean_{diff} = 11$, 95% HDI [1, 21]) has a 98.20% probability of being positive ($> 0$)

and 98.15% of being significant ($> 0.05$). While having "self" slightly increased the RT

than having "other" inside the shapes for the bad character: $mean_{diff} = 5$, 95% HDI [-17,

27]) has a 69.45% probability of being positive ($> 0$) and 69.27% of being significant ($>$

0.05), See Figure 3.

## Discussion

In this study, we investigated the primacy of morality in cognitive processes through

systematically manipulating the factors that are central to the information processing of

morality. First, we found a robust prioritization of good character in response times and $d$'

scores for the shape-label matching tasks across experiments. Second, to pinpoint the

underlying processes of the effect, the analyses revealed that a self-referencing process was

the fundamental driver of these effects, consistent with the self-binding account; that is,

when a stimulus refers to the self, activation of self-representation enhances the binding of

external input with internal knowledge through which self-related information can be

integrated and optimized. The valence account, on the other hand, which posits that the

*Figure 3.* Implicit binding between self and good characters. (A) Posterior distributions of RT under different conditions of Experiment 4a; (B) Posterior distributions of *d*-prime under different conditions of Experiment 4a; (C) Posterior distributions of RT under different conditions of Experiment 4b; (D) Posterior distributions of *d*-prime under different conditions of Experiment 4b; (E) Posterior distributions of the RT differences between good character and neutral character when self (left) and other (right) were presented inside the shapes; (F) Posterior distributions of the *d*-prime differences between good character and neutral character when self (left) and other (right) were presented inside the shapes; (G) Posterior distributions of the RT differences between self- and other-referencing conditions when good character (left) and neutral character (right) were presented inside the shapes; (H) Posterior distributions of the *d*-prime differences between self- and other-referencing conditions when good character (left) and neutral character (right) were presented inside the shapes.

prioritization effect was derived from a general positivity bias towards all (self and others), was not supported by the findings. Importantly, the prioritization effects emerged regardless of whether the relationship between moral character and oneself was task relevant. Collectively, participants tend to attribute moral character to themselves rather than others, leading to prioritized responses to self perceived moral character in decision making.

The current study provided robust evidence for the prioritization of good character in perceptual decision-making. Though the primacy of morality has been argued in social psychology, whether morality is prioritized in information processing has been disputed. For instance, E. Anderson et al. (2011) reported that faces associated with bad social behavior capture attention more rapidly, but an independent team failed to replicate the effect (Stein et al., 2017). In another study, Gantman and Van Bavel (2014) found that moral words are more likely to be judged as words when it was presented subliminally. But this effect may be caused by semantic priming instead of morality (Firestone & Scholl, 2015; Jussim et al., 2016). To overcome this issue, we employed a shape-label matching task to eliminate the semantic priming effect for two reasons. First, associations between shapes and moral characters were acquired during the instruction phase, semantic priming from pre-existed knowledge was impossible (Lee, Martin, & Sui, 2021). Second, there were only a few pairs of stimuli that were used and each stimulus represented different conditions, making it impossible for priming between trials. Importantly, a series of control experiments (1b, 1c, and 2) excluded other confounding factors such as familiarity, presenting sequence, or words-based associations, suggesting that it was the moral content that drove the perceived prioritization of good character. These results are in line with a growing literature on the social and relational nature of perception (Hafri & Firestone, 2021; Xiao, Coppin, & Bavel, 2016).

The prioritization of good character found in the current study was incongruent with previous moral perception studies, which typically reported a negativity bias, i.e.,

information related to bad character is processed preferentially (E. Anderson et al., 2011; Eiserbeck & Abdel Rahman, 2020). This discrepancy may result from different task types employed: while in many moral perception studies, the participants were asked to detect the existence of a stimulus, the current task asked participants to judge the associations between a shape and a person. In other words, previous studies targeted the early stages of perception, while the current task focused more on perceptual decision-making, consistent with previous work (Sui & Humphreys, 2013). This discrepancy is consistent with the positivity bias in studies with emotional stimuli (Pool, Brosch, Delplanque, & Sander, 2016).

The current study expanded previous moral perception studies by testing a novel account that self-referencing processing is the critical driver of the effects. Our results revealed that prioritization of good character is modulated by self-relevance: good character was prioritized when it was referred to oneself. In contrast, good character information was not prioritized when it was referred to others. The modulation effect of self-relevance was amplified when the relationship between moral character and oneself was explicit, consistent with previous studies that only positive aspects of the self are prioritized (Hu et al., 2020). More importantly, the effect persisted even when the relationship between moral character and oneself was task-irrelevant, indicating an implicit self-referencing process emerged from presenting good character and self-related information in the same display. A possible explanation for this spontaneous self-referencing of good character is that the positive moral self-view is central to our identity (Freitas, Cikara, Grossmann, & Schlegel, 2017; Strohminger, Knobe, & Newman, 2017) and the motivation to maintain a moral self-view influences how we perceive (e.g., Ma & Han, 2010) and remember (e.g., Carlson, Maréchal, Oud, Fehr, & Crockett, 2020; Stanley, Henne, & De Brigard, 2019), with implications for the quality of life and wellbeing.

Although the results here revealed the prioritization of good character in perceptual decision-making, we did not claim that the motivation of a moral self-view *penetrates*

perception. The perceptual decision-making process involves processes more than just

encoding the sensory inputs (Scheller & Sui, 2022). To fully account for the nuance of

behavioral data and/or related data collected from other modules (e.g., Sui, He,

Golubickis, Svensson, & Neil Macrae, 2023), we may need computational models and an

integrative experimental approach (Almaatouq et al., 2022). For example, sequential

sampling models suggest that, when making a perceptual decision, the agent continuously

accumulates evidence until the amount of evidence passes a threshold, and then a decision

is made (Chuan-Peng et al., 2022; Forstmann, Ratcliff, & Wagenmakers, 2016; Ratcliff,

Smith, Brown, & McKoon, 2016). In these models, the evidence, or decision variable, can

accumulate from both sensory information but also memory (Shadlen & Shohamy, 2016).

Recently, applications of sequential sample models to perceptual matching tasks also

suggest that different processes may contribute to the prioritization effect of self

(Golubickis et al., 2017) or good self (Hu et al., 2020). Similarly, reinforcement learning

models revealed that the initial discrimination between self- and other-referencing learning

lies in the learning rate (Lockwood et al., 2018). These investigations suggest that

computational models are required to disentangle the cognitive processes underlying the

prioritization of good character.

## References

Abele, A. E., & Bruckmüller, S. (2011). The bigger one of the "Big Two"? Preferential processing of communal information. *Journal of Experimental Social Psychology, 47*(5), 935–948. https://doi.org/10.1016/j.jesp.2011.03.028

Almaatouq, A., Griffiths, T. L., Suchow, J. W., Whiting, M. E., Evans, J., & Watts, D. J. (2022). Beyond Playing 20 Questions with Nature: Integrative Experiment Design in the Social and Behavioral Sciences. *Behavioral and Brain Sciences*, 1–55. https://doi.org/10.1017/S0140525X22002874

Anderson, B. A., Laurent, P. A., & Yantis, S. (2011). Value-driven attentional

550  capture. *Proceedings of the National Academy of Sciences*, *108*(25),

551  10367–10371. https://doi.org/10.1073/pnas.1104047108

552  Anderson, E., Siegel, E. H., Bliss-Moreau, E., & Barrett, L. F. (2011). The visual

553  impact of gossip. *Science*, *332*(6036), 1446–1448.

554  https://doi.org/10.1126/science.1201574

555  Atlas, L. Y. (2023). How Instructions, Learning, and Expectations Shape Pain and

556  Neurobiological Responses. *Annual Review of Neuroscience*, *46*(1).

557  https://doi.org/10.1146/annurev-neuro-101822-122427

558  Bortolon, C., & Raffard, S. (2018). Self-face advantage over familiar and unfamiliar

559  faces: A three-level meta-analytic approach. *Psychonomic Bulletin & Review*,

560  *25*(4), 1287–1300. https://doi.org/10.3758/s13423-018-1487-9

561  Brambilla, M., Sacchi, S., Rusconi, P., & Goodwin, G. P. (2021). The primacy of

562  morality in impression development: Theory, research, and future directions. In

563  B. Gawronski (Ed.), *Advances in Experimental Social Psychology* (Vol. 64, pp.

564  187–262). Academic Press. https://doi.org/10.1016/bs.aesp.2021.03.001

565  Bürkner, P.-C. (2017). Brms: An r package for bayesian multilevel models using

566  stan. *Journal of Statistical Software; Vol 1, Issue 1 (2017)*. Retrieved from

567  https://www.jstatsoft.org/v080/i01 http://dx.doi.org/10.18637/jss.v080.i01

568  Carlson, R. W., Maréchal, M. A., Oud, B., Fehr, E., & Crockett, M. J. (2020).

569  Motivated misremembering of selfish decisions. *Nature Communications*, *11*(1),

570  2100. https://doi.org/10.1038/s41467-020-15602-4

571  Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M.,

572  … Riddell, A. (2017). Stan: A probabilistic programming language [Journal

573  Article]. *Journal of Statistical Software*, *76*(1).

574  https://doi.org/10.18637/jss.v076.i01

575  Carruthers, P. (2021). On valence: Imperative or representation of value? *The

576  British Journal for the Philosophy of Science.* https://doi.org/10.1086/714985

Chuan-Peng, H., Geng, H., Zhang, L., Fengler, A., Frank, M., & Zhang, R.-Y. (2022). *A Hitchhiker's Guide to Bayesian Hierarchical Drift-Diffusion Modeling with dockerHDDM*. PsyArXiv. https://doi.org/10.31234/osf.io/6uzga

Cole, M. W., Braver, T. S., & Meiran, N. (2017). The task novelty paradox: Flexible control of inflexible neural pathways during rapid instructed task learning. *Neuroscience & Biobehavioral Reviews*, *81*, 4–15. https://doi.org/10.1016/j.neubiorev.2017.02.009

Deltomme, B., Mertens, G., Tibboel, H., & Braem, S. (2018). Instructed fear stimuli bias visual attention. *Acta Psychologica*, *184*, 31–38. https://doi.org/10.1016/j.actpsy.2017.08.010

Eiserbeck, A., & Abdel Rahman, R. (2020). Visual consciousness of faces in the attentional blink: Knowledge-based effects of trustworthiness dominate over appearance-based impressions. *Consciousness and Cognition*, *83*, 102977. https://doi.org/10.1016/j.concog.2020.102977

Firestone, C., & Scholl, B. J. (2015). Enhanced visual awareness for morality and pajamas? Perception vs. Memory in "top-down" effects. *Cognition*, *136*, 409–416. https://doi.org/10.1016/j.cognition.2014.10.014

Firestone, C., & Scholl, B. J. (2016a). Cognition does not affect perception: Evaluating the evidence for "top-down" effects. *Behavioral and Brain Sciences*, *39*, e229. https://doi.org/10.1017/S0140525X15000965

Firestone, C., & Scholl, B. J. (2016b). "Moral Perception" Reflects Neither Morality Nor Perception. *Trends in Cognitive Sciences*, *20*(2), 75–76. https://doi.org/10.1016/j.tics.2015.10.006

Fiske, S. T. (1980). Attention and weight in person perception: The impact of negative and extreme behavior. *Journal of Personality and Social Psychology*, *38*(6), 889–906. https://doi.org/10.1037/0022-3514.38.6.889

Forstmann, B. U., Ratcliff, R., & Wagenmakers, E.-J. (2016). Sequential Sampling

Models in Cognitive Neuroscience: Advantages, Applications, and Extensions. *Annual Review of Psychology*, *67*(1). https://doi.org/10.1146/annurev-psych-122414-033645

Freitas, J. D., Cikara, M., Grossmann, I., & Schlegel, R. (2017). Origins of the belief in good true selves. *Trends in Cognitive Sciences*, *21*(9), 634–636. https://doi.org/10.1016/j.tics.2017.05.009

Gantman, A. P., & Bavel, J. J. V. (2015). Moral Perception. *Trends in Cognitive Sciences*, *19*(11), 631–633. https://doi.org/10.1016/j.tics.2015.08.004

Gantman, A. P., & Bavel, J. J. V. (2016). See for Yourself: Perception Is Attuned to Morality. *Trends in Cognitive Sciences*, *20*(2), 76–77. https://doi.org/10.1016/j.tics.2015.12.001

Gantman, A. P., & Van Bavel, J. J. (2014). The moral pop-out effect: Enhanced perceptual awareness of morally relevant stimuli. *Cognition*, *132*(1), 22–29. https://doi.org/10.1016/j.cognition.2014.02.007

Goh, J. X., Hall, J. A., & Rosenthal, R. (2016). Mini meta-analysis of your own studies: Some arguments on why and a primer on how. *Social and Personality Psychology Compass*, *10*(10), 535–549. https://doi.org/10.1111/spc3.12267

Golubickis, M., Falben, J. K., Sahraie, A., Visokomogilski, A., Cunningham, W. A., Sui, J., & Macrae, C. N. (2017). Self-prioritization and perceptual matching: The effects of temporal construal. *Memory & Cognition*, *45*(7), 1223–1239. https://doi.org/10.3758/s13421-017-0722-3

Hafri, A., & Firestone, C. (2021). The Perception of Relations. *Trends in Cognitive Sciences*, *25*(6), 475–492. https://doi.org/10.1016/j.tics.2021.01.006

Haidt, J., & Kesebir, S. (2010). Morality. In S. T. Fiske, D. T. Gilbert, & G. Lindzey (Eds.), *Handbook of Social Psychology* (5th ed., pp. 797–832). John Wiley & Sons, Inc.

Hu, C.-P., Lan, Y., Macrae, C. N., & Sui, J. (2020). Good me bad me: Does valence

influence self-prioritization during perceptual decision-making? *Collabra: Psychology*, *6*(1), 20. https://doi.org/10.1525/collabra.301

Juechems, K., & Summerfield, C. (2019). Where does value come from? *Trends in Cognitive Sciences*, *23*(10), 836–850. https://doi.org/10.1016/j.tics.2019.07.012

Jussim, L., Crawford, J. T., Anglin, S. M., Stevens, S. T., & Duarte, J. L. (2016). Interpretations and methods: Towards a more effectively self-correcting social psychology. *Journal of Experimental Social Psychology*, *66*, 116–133. https://doi.org/10.1016/j.jesp.2015.10.003

Larsen, R. J., Mercer, K. A., & Balota, D. A. (2006). Lexical characteristics of words used in emotional Stroop experiments. *Emotion*, *6*(1), 62–72. https://doi.org/10.1037/1528-3542.6.1.62

Lee, N. A., Martin, D., & Sui, J. (2021). A pre-existing self-referential anchor is not necessary for self-prioritisation. *Acta Psychologica*, *219*, 103362. https://doi.org/10.1016/j.actpsy.2021.103362

Liu, Q., Wang, F., Yan, W., Peng, K., Sui, J., & Hu, C.-P. (2020). Questionnaire data from the revision of a chinese version of free will and determinism plus scale. *Journal of Open Psychology Data*, *8*(1), 1. https://doi.org/10.5334/jopd.49/

Lockwood, P. L., Wittmann, M. K., Apps, M. A. J., Klein-Flügge, M. C., Crockett, M. J., Humphreys, G. W., & Rushworth, M. F. S. (2018). Neural mechanisms for learning self and other ownership. *Nature Communications*, *9*(1), 4747. https://doi.org/10.1038/s41467-018-07231-9

Ma, Y., & Han, S. (2010). Why we respond faster to the self than to others? An implicit positive association theory of self-advantage during implicit face recognition. *Journal of Experimental Psychology: Human Perception and Performance*, *36*, 619–633. https://doi.org/10.1037/a0015797

Makowski, D., Ben-Shachar, M. S., Chen, S. H. A., & Lüdecke, D. (2019). Indices of Effect Existence and Significance in the Bayesian Framework. *Frontiers in*

658    *Psychology, 10.* https://doi.org/10.3389/fpsyg.2019.02767

659    Makowski, D., Ben-Shachar, M. S., & Lüdecke, D. (2019). bayestestR: Describing

660        Effects and their Uncertainty, Existence and Significance within the Bayesian

661        Framework. *Journal of Open Source Software, 4*(40), 1541.

662        https://doi.org/10.21105/joss.01541

663    Ohman, A., Lundqvist, D., & Esteves, F. (2001). The face in the crowd revisited: A

664        threat advantage with schematic stimuli. *Journal of Personality and Social*

665        *Psychology, 80*(3), 381–396. https://doi.org/10.1037/0022-3514.80.3.381

666    Pool, E., Brosch, T., Delplanque, S., & Sander, D. (2016). Attentional bias for

667        positive emotional stimuli: A meta-analytic investigation. *Psychological*

668        *Bulletin, 142*(1), 79–106. https://doi.org/10.1037/bul0000026

669    Ratcliff, R., Smith, P. L., Brown, S. D., & McKoon, G. (2016). Diffusion Decision

670        Model: Current Issues and History. *Trends in Cognitive Sciences, 20*(4),

671        260–281. https://doi.org/10.1016/j.tics.2016.01.007

672    Rouder, J. N., & Lu, J. (2005). An introduction to bayesian hierarchical models

673        with an application in the theory of signal detection. *Psychonomic Bulletin &*

674        *Review, 12*(4), 573–604. https://doi.org/10.3758/bf03196750

675    Rousselet, G. A., & Wilcox, R. R. (2020). Reaction times and other skewed

676        distributions: Problems with the mean and the median. *Meta-Psychology, 4.*

677        https://doi.org/10.15626/MP.2019.1630

678    Scheller, M., & Sui, J. (2022). The power of the self: Anchoring information

679        processing across contexts. *Journal of Experimental Psychology: Human*

680        *Perception and Performance, 48*(9), 1001–1021.

681        https://doi.org/10.1037/xhp0001017

682    Shadlen, M. N., & Shohamy, D. (2016). Decision Making and Sequential Sampling

683        from Memory. *Neuron, 90*(5), 927–939.

684        https://doi.org/10.1016/j.neuron.2016.04.036

Shore, D. M., & Heerey, E. A. (2013). Do social utility judgments influence attentional processing? *Cognition*, *129*(1), 114–122. https://doi.org/10.1016/j.cognition.2013.06.011

Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2013). *Life after p-hacking* [Conference Proceedings]. https://doi.org/10.2139/ssrn.2205186

Stanley, M. L., Henne, P., & De Brigard, F. (2019). Remembering moral and immoral actions in constructing the self. *Memory & Cognition*, *47*(3), 441–454. https://doi.org/10.3758/s13421-018-0880-y

Stein, T., Grubb, C., Bertrand, M., Suh, S. M., & Verosky, S. C. (2017). No impact of affective person knowledge on visual awareness: Evidence from binocular rivalry and continuous flash suppression. *Emotion*, *17*(8), 1199–1207. https://doi.org/10.1037/emo0000305

Strohminger, N., Knobe, J., & Newman, G. (2017). The True Self: A Psychological Concept Distinct From the Self. *Perspectives on Psychological Science*, *12*(4), 551–560. https://doi.org/10.1177/1745691616689495

Sui, J., He, X., Golubickis, M., Svensson, S. L., & Neil Macrae, C. (2023). Electrophysiological correlates of self-prioritization. *Consciousness and Cognition*, *108*, 103475. https://doi.org/10.1016/j.concog.2023.103475

Sui, J., He, X., & Humphreys, G. W. (2012). Perceptual effects of social salience: Evidence from self-prioritization effects on perceptual matching. *Journal of Experimental Psychology: Human Perception and Performance*, *38*(5), 1105–1117. https://doi.org/10.1037/a0029792

Sui, J., & Humphreys, G. W. (2013). The boundaries of self face perception: Response time distributions, perceptual categories, and decision weighting. *Visual Cognition*, *21*(4), 415–445. https://doi.org/10.1080/13506285.2013.800621

Sui, J., & Humphreys, G. W. (2015a). More of me! Distinguishing self and reward bias using redundancy gains. *Attention, Perception, & Psychophysics*, *77*(8),

2549–2561. https://doi.org/10.3758/s13414-015-0970-x

Sui, J., & Humphreys, G. W. (2015b). The Integrative Self: How Self-Reference

Integrates Perception and Memory. *Trends in Cognitive Sciences*, *19*(12),

719–728. https://doi.org/10.1016/j.tics.2015.08.015

Sui, J., & Rotshtein, P. (2019). Self-prioritization and the attentional systems.

*Current Opinion in Psychology*, *29*, 148–152.

https://doi.org/10.1016/j.copsyc.2019.02.010

Unkelbach, C., Alves, H., & Koch, A. (2020). Chapter three - negativity bias,

positivity bias, and valence asymmetries: Explaining the differential processing

of positive and negative information. In B. Gawronski (Ed.), *Advances in*

*experimental social psychology* (Vol. 62, pp. 115–187). Academic Press.

https://doi.org/10.1016/bs.aesp.2020.04.005

Xiao, Y. J., Coppin, G., & Bavel, J. J. V. (2016). Perceiving the world through

group-colored glasses: A perceptual model of intergroup relations. *Psychological*

*Inquiry*, *27*(4), 255–274. https://doi.org/10.1080/1047840X.2016.1199221

Yaoi, K., Osaka, M., & Osaka, N. (2021). Does Implicit Self-Reference Effect Occur

by the Instantaneous Own-Name? *Frontiers in Psychology*, *12*, 4440.

https://doi.org/10.3389/fpsyg.2021.709601

Ybarra, O., Chan, E., & Park, D. (2001). Young and old adults' concerns about

morality and competence. *Motivation and Emotion*, *25*, 85–100.

https://doi.org/10.1023/A:1010633908298