1    Priorization of the morally good depends on self-relevance in perceptual matching

2    Hu Chuan-Peng[1,2], Kaiping Peng[2], & Jie Sui[3]

3    [1] Nanjing Normal University, 210024 Nanjing, China

4    [2] Tsinghua University, 100084 Beijing, China

5    [3] University of Aberdeen, Aberdeen, Scotland

6    Author Note

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

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

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

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

11 HCP analyzed the data and drafted the manuscript. All authors read and agreed upon the

12 current version of the manuscripts.

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

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

15 210024 Nanjing, China. E-mail: hcp4715@gmail.com

16                                    Abstract

17   To navigate in a complex social world, our cognitive system are evolved to be sensitive to

18   social information. Among all these social informaiton, morality related information is of

19   special interest. On the one hand, paying attention to other's moral character profitable for

20   ourselves. On the other hand, we need to maitain a moral self-view that fit the soical

21   norm. Though behavioral effects of moral character and moral self-enhancement had been

22   extensively studied in psychology of morality, social perception, and identity, whether the

23   moral character related information can impact low-level perceptual process is unknown.

24   In a series of experiments, we examined the effect of immediately acquired moral character

25   information on perceptual matching. Participants first learned the association between

26   moral character and visual cues (shapes), then performed a perceptual matching task. The

27   results showed that shapes associated with positive moral character were prioritized, as

28   compared to neutral or negative bad moral characters. This pattern was robust after

29   changing the words for moral charachter or using diagnostic behavioral as an proxy of

30   mroal character. Also, this patterns were robust when changing simultaneous presentation

31   to sequential presentation. We then examined two approximate explanations for this effect:

32   value-based prioritization or social-categorization based prioritization. We manipulated the

33   identity of different moral character explicitly and found that the good moral character

34   effect was strong when for the self-referential conditions but weak or non-exist for

35   other-referential condition. We further tested the good-self based social categorization by

36   presenting the identity or moral character information as task-irrelevant stimuli, so that we

37   can distinguish between the unique good-self hypothesis and a more general good-person

38   based social categorization hypothesis. We found that ….., these results suggested that

39   participants are more senstive to the moral valence of self when the valence were

40   task-irrelevant, but less sensitive to the identity of the morally good when the identity were

41   task-irrelvant. These results added new evidence for the social vision and suggested the

42   advantage of moral good depends on the self-relevant in perceptual decision-making task,

43    instead of perspective free.

44            *Keywords:* Perceptual decision-making, Self positivity bias, moral character

45            Word count: X

<sup>46</sup> Priorization of the morally good depends on self-relevance in perceptual matching

## Introduction

<sup>48</sup> [sentences in bracket are key ideas]

<sup>49</sup> social vision –> moral vision –> two competing explanations (value-based
<sup>50</sup> vs. true-self-based) –> true-self is not perspective free but self-centered.

<sup>51</sup> Will not include experiment 5; stop exploring the correlations.

<sup>52</sup> Our information processing system had been evolved in a way that top-down factor
<sup>53</sup> can modulate the low-level processes. There are debates on whether perception can be
<sup>54</sup> penetrated by top-down factors.

<sup>55</sup> [Morality is the central of human social life]. People experience a substantial amount
<sup>56</sup> of moral events in everyday life (e.g., Hofmann, Wisneski, Brandt, & Skitka, 2014). When
<sup>57</sup> experiencing these events, it always involves judging "good" or "bad." By judging "right"
<sup>58</sup> or "wrong," people are implicitly judging the moral character of involved parties as "good"
<sup>59</sup> vs. "bad" (Uhlmann, Pizarro, & Diermeier, 2015). The central role of moral character also
<sup>60</sup> supported by the extensive studies from person perception and social evaluation, where
<sup>61</sup> morality is a basic dimension for social evaluation (Abele, Ellemers, Fiske, Koch, &
<sup>62</sup> Yzerbyt, 2020; Goodwin, 2015; Goodwin, Piazza, & Rozin, 2014; Willis & Todorov, 2006)
<sup>63</sup> and the most important aspect to evaluate the continuity of identity (Strohminger, Knobe,
<sup>64</sup> & Newman, 2017).

<sup>65</sup> Given the importance of moral character, to successfully navigate in a social world, a
<sup>66</sup> person needs to both evaluate others' moral character and behave in a way that she/he is
<sup>67</sup> perceived as a moral person, or at least not a morally bad person. Maintaining a moral
<sup>68</sup> self-view is as important as making judgments about others' moral character (Ellemers,
<sup>69</sup> Toorn, Paunov, & Leeuwen, 2019). Indeed, previous studies found that people maintain a
<sup>70</sup> positive moral self-view even after dishonest behavior (Monin & Jordan, 2009) and that

people evaluate themselves as morally superior to others (Klein & Epley, 2016; Tappin &
McKay, 2017). Recent theorists further integrated the moral judgment and moral self-view,
proposed a person-centered account for moral psychology, which focused on the individuals
in moral evaluation instead of acts (Uhlmann, Pizarro, & Diermeier, 2015). Under this
framework, previous seemingly contradicting phenomenons can be explained. For example,
whether people decide to expose an unethical behavior depends on how their relationship
of the target (e.g., Waytz, Dungan, & Young, 2013).

To date, however, as Freeman and Ambady (2011) put it, studies in the perception of
moral character didn't try to explain the perceptual process, rather, they are trying to
explain the higher-order social cognitive processes that come after. Essentially, these
studies are perception of moral character without perceptual process. Without knowledge
of perceptual processes, we can not have a full picture of how moral character is processed
in our cognition. As an increasing attention is paid to perceptual process underlying social
cognition, it's clear that perceptual processes are strongly influenced by social factors, such
as group-categorization, stereotype (see Bagnis, Celeghin, Mosso, & Tamietto, 2019; Stolier
& Freeman, 2016; Xiao, Coppin, & Bavel, 2016). Given the importance of moral character
and that moral character related information has strong influence on learning and memory
(Carlson, Maréchal, Oud, Fehr, & Crockett, 2020; Stanley & De Brigard, 2019), one might
expect that moral character related information could also play a role in perceptual process.

To explore the perceptual process of moral character and the underlying mechanism,
we conducted a series of experiments to explore (1) whether we can detect the influence of
moral character information on perceptual decision-making in a reliable way, and (2)
potential explanations for the effect. In the first four experiments, we found a robust effect
of good-person prioritization in perceptual decision-making. Then, we explore the potential
explanations and tested value-based prioritization versus good-self based prioritization
(social-categorization (Turner, Hogg, Oakes, Reicher, & Wetherell, 1987; Turner, Oakes,
Haslam, & McGarty, 1994)). These results suggested that people may categorize self and

⁹⁸ other based on moral character; in these categorizations, the core self, i.e., the good-self, is
⁹⁹ always prioritized.

**Perceptual process of moral character**

¹⁰¹     [exp1a, b, c, and exp2]

¹⁰² [using associative learning task to study the moral character's influence on
¹⁰³ perception] Though it is theoretically possible that moral character related information
¹⁰⁴ may be prioritized in perceptual process, no empirical studies had directly explored this
¹⁰⁵ possibility. One difficulty of studying the perceptual process of moral character is that
¹⁰⁶ moral character is an inferred trait instead of observable feature. Usually, one needs more
¹⁰⁷ sensory input, e.g., behavior history, to infer moral character of a person. For example,
¹⁰⁸ Anderson, Siegel, Bliss-Moreau, and Barrett (2011) asked participant to first study the
¹⁰⁹ behavioral description of faces and then asked them to perform a perceptual detection task.
¹¹⁰ They assumed that by learning the behavioral description of a person (represented by a
¹¹¹ face), participants can acquire the moral related information about faces, and the
¹¹² associations could then bias the perceptual processing of the faces (but see Stein, Grubb,
¹¹³ Bertrand, Suh, and Verosky (2017)). One drawback of this approach is that participants
¹¹⁴ may differ greatly when inferring the moral character of the person from behavioral
¹¹⁵ descriptions, given that notion what is morality itself is varying across population Jones et
¹¹⁶ al. (2021) and those descriptions and faces may themselves are idiosyncratic, therefore,
¹¹⁷ introduced additional variance to the targeted effect.

¹¹⁸     An alternative is to use abstract semantic concepts. Abstract concepts of moral
¹¹⁹ character are used to describe and represent moral characters. These abstract concepts
¹²⁰ may be part of a dynamic network in which sensory cue, concrete behaviors and other
¹²¹ information can activate/inhibit each other (e.g., aggressiveness) (Amodio, 2019; Freeman
¹²² & Ambady, 2011). If a concept of moral character (e.g., good person) is activated, it

should be able to influence on the perceptual process of the visual cues through the

dynamic network, especially when the perceptual decision-making is about the concept-cue

association. In this case, abstract concepts of moral character may serve as signal of moral

reputation (for others) or moral self-concept. Indeed, previous studies used the moral

words and found that moral related information can be perceived faster Firestone & Scholl

(2015). If moral character is an important in person perception, then, just as those other

information such as races and stereotype (see Xiao, Coppin, & Bavel, 2016), moral

character related concepts also change the perceptual processes.

To investigate the above possibility, we used an associative learning paradigm to

study how moral character concept change perceptual decision-making. In this paradigm,

simple geometric shapes were paired with different words whose dominant meaning is

describing the moral character of a person. Participants first learn the associations between

shapes and words, e.g., triangle is a good-person. After formed direct associations between

the labels of moral characters and visual cues, participants performed a perceptual

matching task to judge whether the shape-word pair presented on the screen match the

association they learned. This paradigm has been used in studying the perceptual process

of self-concept, but had also proven useful in studying other concepts like social group

(e.g., Enock, Hewstone, Lockwood, & Sui, 2020). By using simple and morally neutral

shapes, we controlled the variations caused by visual cues.

Our first question is, whether the words used the in the associative paradigm is really

related to the moral character? This assumption is consistent with previous theories,

especially the interactive dynamic theory. To validate that moral character concepts

activated moral character as a social cue, we used four experiments to explore and validate

the paradigm. The first experiment directly adopted associative paradigm and changed

labels from "self," "friend," and "stranger" to "good-person," "neutral-person," and

"bad-person." We further tried semantic labels that have more explicit moral meaning

("kind-person," "neutral-person," and "evil-person"). In the third experiments, as in

Anderson, Siegel, Bliss-Moreau, and Barrett (2011), we asked participant to learn the association between three different diagnostic behavior and three different names, and then use the names as moral labels for the associative learning. Finally, we also tested that simultaneously present shape-word pair and sequentially present word and shape didn't change the pattern. All of these four experiments showed a consistent pattern of effect, that is, the visual cues that associated with positive moral character were prioritized.

## Morality as a social-categorization?

[possible explanations: person-based self-categorization vs. stimuli-based valence] The robust pattern from the first four experiments revealed a novel pattern that needs an explanation. It's novel because it's contradict with the "negative is stronger than positive" hypothesis in social psychology (Baumeister, Bratslavsky, Finkenauer, & Vohs, 2001). There are two major alternatives. One possible explanation is the value-based attention, which suggested that valuable stimuli is prioritized in our low-level cognitive processes. Because positive moral character is potentially rewarding, e.g., potential cooperators, it is valuable to individuals and therefore being prioritized. Most empirical evidence for value-based attention are from experiments used monetary reward. However, the monetary reward might be different greatly from the morality in social setting. So far, only a few empirical studies supported the value-based attention in social evaluation. Eiserbeck and Abdel Rahman (2020) and Shore and Heerey (2013) found that neutral faces, after associated with trustworthy behavioral description or trustworthy interaction history, attracted attention more than untrustworthy faces, probably because trustworthy faces are more likely to be the collaborative partners subsequent tasks, which will bring reward. Applying this explanation to the current setting need a further assumption that participants automatically view the moral character related information as self-relevant objects. Only based on the objectified stimuli that we evaluate their value (rewarding or threatening) to us (Juechems & Summerfield, 2019; Reicher & Hopkins, 2016).

Another possibility is that we will perceive those moral character not as objects but as person, and automatic categorize whether they are in-group or out-group, instead of calculating their value to us. This account assumed that moral character served as a way to categorize other. In the first four experiments' situation, the identity of the moral character is ambiguous, participants may automatically categorize morally good people as in-group (as an extension of themselves) and therefore preferentially processed these information.

However, the above four experiments could not distinguish between these two possibilities, because the concept "good-person" can both be rewarding and be categorized as in-group member, and previous studies using associative learning paradigm revealed that both rewarding stimuli (e.g., Sui, He, & Humphreys, 2012) and in-group information (Enock, Hewstone, Lockwood, & Sui, 2020) are prioritized.

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

203     We introduced identity (self vs. other) as an addition independent variable in exp 3a,

204 3b, and 6b. Now the moral valence is orthogonal to the identity. We found that (1)

205 good-self is always faster than neutral-self and bad-self, but good-other only have weak to

206 null advantage to neutral-other and bad-other. which mean the social categorization is

207 self-centered. (2) good-self's advantage over good other only occur when self- and other-

208 were in the same task. i.e. the relative advantage is competition based instead of absolute.

209 These three experiments suggest that people more like to view the moral character stimuli

210 as person and categorize good-self as an unique category against all others. A three-level

211 Bayesian generalized linear mixed effect model showed that there was no effect of valence

212 when the identity was other. This results showed that value-based attention was not likely

213 the mechanism behind the pattern we observed in first four experiments. However, it is

214 still unclear Why good-self was prioritized. Besides the social-categorization explanation,

215 it's also possible that good self is so unique that it is prioritized in all possible situation

216 and therefore is not social categorization *per se.*

217     [what we care? valence of the self exp4a or identity of the good exp4b?] We go

218 further to disentangle the good-self complex: is it because the special role of good-self or

219 because of social categorization. We designed two complementary experiments. in

220 experiment 4a, participants only learned the association between self and other, the words

221 "good-person," "neutral person," and "bad person" were presented as task-irrelevant

222 stimuli, while in experiment 4b, participants learned the associations between

223 "good-person," "neutral-person," and "bad-person," and the "self" and "other" were

224 presented as task-irrelevant stimuli. These two experiment can be used to distinguish the

225 "good-self" as anchor account and the "good-self-based social categorization" account. If

226 good-self as an anchor is true, then, in both experiment, good-self will show advantage over

227 all other stimuli. More specifically, in experiment 4a, where only the self-relevance is

228 task-relevant, there will be advantage for good as task-irrelevant condition than the other

229 two self conditions; in experiment 4b, in the good condition, there will be an advantage for

self as task-irrelevant condition over other as task-irrelevant condition. If good-self-based social categorization if true, then, the prioritization effect will depends on whether the stimuli can be categorized as the same group of good-self. More specifically, in experiment 4a, there will be good-as-task-irrelevant stimuli than other condition in self conditions, this prediction is the same as the "good-self as anchor" account; however, for experiment 4b, there will be no self-as-task-irrelevant stimuli than other-as-task-irrelevant condition.

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

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

Key concepts and discussing points:

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

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

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

**Variable self**: Who we are, how we see ourselves, how we define our relations to others (indeed whether they are construed as 'other' or as part of the extended 'we' self) is

different in different settings.

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

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

What is the relation between morally good and self in a semantic network (attractor network) (Freeman & Ambady, 2011)? The psychological essentialism account proposed that the moral good self is perspective independent, i.e., there is a moral good self in all. This perspective free effect is not exist in our effect.

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

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

in experiment 4a and 4b, we didn't have a baseline condition where there is no word inside the shape?

## Disclosures

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

All the experiments reported were not pre-registered. Most experiments (1a ~ 6b, except experiment 3b) reported in the current study were first finished between 2014 to 2016 in Tsinghua University, Beijing, China. Participants in these experiments were

280 recruited in the local community. To increase the sample size of experiments to 50 or more

281 (Simmons, Nelson, & Simonsohn, 2013), we recruited additional participants in Wenzhou

282 University, Wenzhou, China in 2017 for experiment 1a, 1b, 4a, and 4b. Experiment 3b was

283 finished in Wenzhou University in 2017. To have a better estimation of the effect size, we

284 included the data from two experiments (experiment 7a, 7b) that were reported in Hu,

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

286 All participant received informed consent and compensated for their time. These

287 experiments were approved by the ethic board in the Department of Tsinghua University.

## General methods

### Design and Procedure

290 This series of experiments studied the perceptual process of moral character, using

291 the social associative learning paradigm (or tagging paradigm)(Sui, He, & Humphreys,

292 2012), in which participants first learned the associations between geometric shapes and

293 labels of person with different moral character (e.g., in first three studies, the triangle,

294 square, and circle and good person, neutral person, and bad person, respectively). The

295 associations of the shapes and label were counterbalanced across participants. After

296 remembered the associations, participants finished a practice phase to familiar with the

297 task, in which they viewed one of the shapes upon the fixation while one of the labels below

298 the fixation and judged whether the shape and the label matched the association they

299 learned. When participants reached 60% or higher accuracy at the end of the practicing

300 session, they started the experimental task which was the same as in the practice phase.

301 The experiment 1a, 1b, 1c, 2, and 6a shared a 2 (matching: match vs. nonmatch) by

302 3 (moral character: good person vs. neutral person vs. bad person) within-subject design.

303 Experiment 1a was the first one of the whole series studies and found the prioritization of

304 stimuli associated with good-person. To confirm that it is the moral character that caused

the effect, we further conducted experiment 1b, 1c, and 2. More specifically, experiment 1b used different Chinese words as label to test whether the effect only occurred with certain familiar words. Experiment 1c manipulated the moral valence indirectly: participants first learned to associate different moral behaviors with different neutral names, after remembered the association, they then performed the perceptual matching task by associating names with different shapes. Experiment 2 further tested whether the way we presented the stimuli influence the effect of valence, by sequentially presenting labels and shapes. Note that part of participants of experiment 2 were from experiment 1a because we originally planned a cross task comparison. Experiment 6a, which shared the same design as experiment 2, was an EEG experiment which aimed at exploring the neural correlates of the effect. But we will focus on the behavioral results of experiment 6a in the current manuscript.

For experiment 3a, 3b, 4a, 4b, 6b, 7a, and 7b, we included self-reference as another within-subject variable in the experimental design. For example, the experiment 3a directly extend the design of experiment 1a into a 2 (matchness: match vs. nonmatch) by 2 (reference: self vs. other) by 3 (moral valence: good vs. neutral vs. bad) within-subject design. Thus in experiment 3a, there were six conditions (good-self, neutral-self, bad-self, good-other, neutral-other, and bad-other) and six shapes (triangle, square, circle, diamond, pentagon, and trapezoids). The experiment 6b was an EEG experiment extended from experiment 3a but presented the label and shape sequentially. Because of the relatively high working memory load (six label-shape pairs), experiment 6b were conducted in two days: the first day participants finished perceptual matching task as a practice, and the second day, they finished the task again while the EEG signals were recorded. Experiment 3b was designed to separate the self-referential trials and other-referential trials. That is, participants finished two different types of block: in the self-referential blocks, they only responded to good-self, neutral-self, and bad-self, with half match trials and half non-match trials; in the other-reference blocks, they only responded to good-other,

neutral-other, and bad-other. Experiment 7a and 7b were designed to test the cross task

robustness of the effect we observed in the aforementioned experiments (see, Hu, Lan,

Macrae, & Sui, 2020). The matching task in these two experiments shared the same design

with experiment 3a, but only with two moral character, i.e., good vs. bad. We didn't

include the neutral condition in experiment 7a and 7b because we found that the neutral

and bad conditions constantly showed non-significant results in experiment $1 \sim 6$.

Experiment 4a and 4b were design to explore the mechanism behind the

prioritization of good-self. In 4a, we used only two labels (self vs. other) and two shapes

(circle, square). To manipulate the moral valence, we added the moral-related words within

the shape and instructed participants to ignore the words in the shape during the task. In

4b, we reversed the role of self-reference and valence in the task: participant learnt three

labels (good-person, neutral-person, and bad-person) and three shapes (circle, square, and

triangle), and the words related to identity, "self" or "other," were presented in the shapes.

As in 4a, participants were told to ignore the words inside the shape during the task.

E-prime 2.0 was used for presenting stimuli and collecting behavioral responses,

except that experiment 7a and 7b used Matlab Psychtoolbox (Brainard, 1997; Pelli, 1997).

For participants recruited in Tsinghua University, they finished the experiment individually

in a dim-lighted chamber, stimuli were presented on 22-inch CRT monitors and their head

were fixed by a chin-rest brace. The distance between participants' eyes and the screen was

about 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 the shape or the word and the fixation cross was

$3.5°$ of visual angle. For participants recruited in Wenzhou University, they finished the

experiment in a group consisted of $3 \sim 12$ participants in a dim-lighted testing room.

Participants were required to finished the whole experiment independently. Also, they were

instructed to start the experiment at the same time, so that the distraction between

participants were minimized. The stimuli were presented on 19-inch CRT monitor. The

359    visual angles are could not be exactly controlled because participants' chin were not fixed.

360    In most of these experiments, participant were also asked to fill a battery of

361    questionnaire after they finish the behavioral tasks. All the questionnaire data are open

362    (see, dataset 4 in Liu et al., 2020). See Table S1 for a summary information about all the

363    experiments.

## Data analysis

365    We used the `tidyverse` of r (see script `Load_save_data.r`) to exclude the practicing

366    trials, invalid trials of each participants, and invalid participants, if there were any, in the

367    raw data. Results of each experiment were then analyzed in two Bayesian approaches and

368    reported in supplementary materials.

369    ***Bayesian hierarchical model.***    We first tested the effect of experimental

370    manipulation using Bayesian hierarchical model. More specifically, we used the Bayesian

371    hierarchical model (BHM, or Bayesian generalized linear mixed model, BGLMM) to model

372    the reaction time and accuracy data. We used Bayesian hierarchical model because BHM

373    provided three advantages over the classic NHST approach (repeated measure ANOVA or

374    $t$-tests): first, BHM estimate the posterior distributions of parameters for statistical

375    inference, therefore provided uncertainty in estimation (Rouder & Lu, 2005). Second,

376    BHM, as generalized linear mixed models, can use distribution that fit the distribution of

377    real data instead of using normal distribution for all data. Using appropriate distributions

378    for the data will avoid misleading results and provide better fitting of the data. For

379    example, Reaction times are not normally distributed but right skewed, and the linear

380    assumption in ANOVAs is not satisfied (Rousselet & Wilcox, 2019). Third, BHM provided

381    an unified framework to analyze data from different levels and different sources, avoid the

382    information loss when we need to combine data from different levels.

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

384 2017) to sample from the posterior, to build the model for RTs and accuracy separately.

385 Using the Bayesian hierarchical model, we can directly estimate the over-all effect across

386 similar experiments with similar experimental design, instead of using a two-step approach

387 where we first estimate parameters, e.g., $d'$ for each participant, and then use a random

388 effect model meta-analysis to synthesize the effect (Goh, Hall, & Rosenthal, 2016). We also

389 we used HDDM to model RTs and accuracy data together using drift diffusion model as

390 the data generative model.

391    *Accuracy.*   We followed practice of previous studies (Hu, Lan, Macrae, & Sui, 2020;

392 Sui, He, & Humphreys, 2012) and used signal detection theory approach to analyze the

393 accuracy data. More specifically, the match trials are treated as signal and the non-match

394 trials are noise. As we mentioned above, we estimated the sensitivity and criterion of SDT

395 by BHM (Rouder & Lu, 2005). Because the BHM can model different level's data using a

396 single unified model, we used a three-level HBM to model the valence effect, which include

397 five experiments: 1a, 1b, 1c, 2, and 6a. Also, we modelled the experiments with both

398 identity and moral valence with a three-level HBM model, which includes 3a, 3b, and 6b.

399 For experiment 4a and 4b, we used two-level models for each separately. However, we

400 compared the posterior of parameters directly because we have full posterior distribution of

401 the effect and can directly compare the posteriors.

402    We used the Bernoulli distribution to model the accuracy data. For a single

403 participant, we assume that the accuracy of $i$th trial is Bernoulli distributed (binomial with

404 1 trial), with probability $p_i$ that $y_i = 1$.

$$y_i \sim Bernoulli(p_i)$$

405 and the probability of choosing "match" $p_i$ at the $i$th trial is a function of the trial type:

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

406  therefore, the outcomes $y_i$ are 0 if the participant responded "nonmatch" on the $i$th trial, 1

407  if they responded "match." We then write the generalized linear model on the probits

408  (z-scores; $\Phi$, "Phi") of $p$s. $\Phi$ is the cumulative normal density function and maps $z$ scores

409  to probabilities. In this way, the intercept of the model ($\beta_0$) is the standardized false alarm

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

411  slope of the model ($\beta_1$) is the increased probability of responding "match" when the trial

412  type is "match," in $z$-scores, which is another expression of $d'$. Therefore, $c = $ -$z$HR $=$

413  $-\beta_0$, and $d' = \beta_1$.

414      In our experimental design, there are three conditions for both match and non-match

415  trials, we can estimate the $d'$ and $c$ separately for each condition. In this case, the criterion

416  $c$ is modeled as the main effect of valence, and the $d'$ can be modeled as the interaction

417  between valence and match, and we explicitly removed the intercept:

$$\Phi(p_i) = 0 + \beta_0 Valence_i + \beta_1 IsMatch_i * Valence_i$$

418      In each experiment, we had multiple participants. We can estimate the group-level

419  parameters by extending the above model into a two-leve model, where we can estimate

420  parameters on individual level and the group level parameter simultaneously. The

421  probability that the $j$th subject responded "match" ($y_{ij} = 1$) at the $i$th trial $p_{ij}$. In the

422  same vein, we have

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

423  The the generalized linear model can be re-written to include two levels:

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

424  We again can write the generalized linear model on the probits (z-scores; $\Phi$, "Phi") of $p$s.

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

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

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

427    For experiments that had 2 (matching: match vs. non-match) by 3 (moral character:

428    good vs. neutral vs. bad), i.e., experiment 1a, 1b, 1c, 2, 5, and 6a, the formula for accuracy

429    in BRMs is as follow:

430    `saymatch ~ 0 + Valence + Valence:ismatch + (0 + Valence +`

431    `Valence:ismatch | Subject), family = bernoulli(link="probit")`

432    For experiments that had two by two by three design, we used the follow formula for

433    the BGLM:

434    `saymatch ~ 0 + ID:Valence + ID:Valence:ismatch + (0 + ID:Valence +`

435    `ID:Valence:ismatch | Subject), family = bernoulli(link="probit")`

436    In the same vein, we can estimate the posterior of parameters across different

437    experiments. We can use a nested hierarchical model to model all the experiment with

438    similar design:

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

439    the generalized linear model is then

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

440    The outcomes $y_{ijk}$ are 0 if participant $j$ in experiment k responded "mismatch" on trial $i$, 1

441    if they responded "match."

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

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

443   distribution:

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

444   in which $mu_0$ and $mu_1$ means the population level parameter.

445     *Reaction times.*   For the reaction time, we used the log normal distribution

446   (https://lindeloev.github.io/shiny-rt/#34_(shifted)_log-normal) to model the data. This

447   means that we need to estimate the posterior of two parameters: $\mu$, $\sigma$. $\mu$ is the mean of the

448   `logNormal` distribution, and $\sigma$ is the disperse of the distribution. Although the log normal

449   distribution can be extended to shifted log normal distribution, with one more parameter:

450   shift, which is the earliest possible response, we found that the additional parameter didnt'

451   improved the model fitting and therefore used the logNormal in our final analysis.

452     The reaction time of the $j$th subject on $i$th trial is a linear function of trial type:

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

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

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

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

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

$$\sigma_j \sim Cauchy()$$

455   Formula used for modeling the data as follow:

```
RT_sec ~ Valence*ismatch + (Valence*ismatch | Subject), family =
lognormal()
```

or

```
RT_sec ~ ID*Valence*ismatch + (ID*Valence*ismatch | Subject), family =
lognormal()
```

we expanded the RT model three-level model in which participants and experiments are two group level variable and participants were nested in the experiments.

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

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

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

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

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

$$\theta_k \sim Cauchy()$$

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

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

474    *Implicit interaction between valence and self-relevance.*    In the third part, we focused

475    on experiment 4a and 4b, which were designed to examine two more nuanced explanation

476    concerning the good-self. The design of experiment 4a and 4b are complementary.

477    Together, they can test whether participants are more sensitive to the moral character of

478    the Self (4a), or the identity of the morally Good (4b).

479    Finally, we explored correlation between results from behavioral results and

480    self-reported measures.

481    For the questionnaire part, we are most interested in the self-rated distance between

482    different person and self-evaluation related questionnaires: self-esteem, moral-self identity,

483    and moral self-image. Other questionnaires (e.g., personality) were not planned to

484    correlated with behavioral data were not included. Note that all questionnaire data were

485    reported in (Liu et al., 2020).

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

487    drift rate ($v$), boundary separation ($a$), and non decision-making time ($t$), because these

488    parameters has relative clear psychological meaning. We used the mean of parameter

489    posterior distribution as the estimate of each parameter for each participants in the

490    correlation analysis. We used alpha = 0.05 and used bootstrap by `BootES` package (Kirby

491    & Gerlanc, 2013) to estimate the correlation.

492    **Hierarchical drift diffusion model (HDDM).**    To further explore the

493    psychological mechanism under perceptual decision-making, we used a generative mode

494    drift diffusion model (DDM) to model our RTs and accuracy data. As the hypothesis

495    testing part, we also used hierarchical Bayesian model to fit the DDM. The package we

496    used was the HDDM (Wiecki, Sofer, & Frank, 2013), a python package for fitting

497    hierarchical DDM. We used the prior implemented in HDDM, that is, weakly informative

498    priors that constrains parameter estimates to be in the range of plausible values based on

499    past literature (Matzke & Wagenmakers, 2009). As reported in Hu, Lan, Macrae, and Sui

500  (2020), we used the stimulus code approach, match response were coded as 1 and

501  nonmatch responses were coded as 0. To fully explore all parameters, we allow all four

502  parameters of DDM free to vary. We then extracted the estimation of all the four

503  parameters for each participants for the correlation analyses. However, because the

504  starting point is only related to response (match vs. non-match) but not the valence of the

505  stimuli, we didn't included it in correlation analysis.

506  **Part 1: Perceptual processing moral character related information**

507       In this part, we report results from five experiments that tested whether an

508  associative learning task, including 192 participants. Note that for both experiment 1a and

509  1b, there were two independent samples with different equipment, trials numbers and

510  testing situation. Therefore, we modeled them as independent samples. These five

511  experiments revealed a robust effect of moral character on perceptual matching task.

512       For the $d$ prime, we found robust effect of moral valence. Shapes associated with

513  positive moral valence ("good person," "kind person" or a name associated with morally

514  good behavioral history) has higher sensitivity (mean = , 95% HDI = ) than shapes

515  associated with neutral condition (mean = , 95% HDI = ), but we did not find differences

516  between shapes associated with negative moral label (mean = , 95% HDI = ) and neutral

517  condition.

518       For the reaction times, we also found robust effect of moral valence. Shapes

519  associated with positive moral valence has faster responses (mean = , 95% HDI = ) than

520  shapes associated with neutral condition (mean = , 95% HDI = ). We also found that the

521  responses to shapes associated with negative moral valence (mean = , 95% HDI = ) were

522  slower as compared to the neutral condition. See Figure 1.

### Part 2: interaction between valence and identity

In this part, we report three experiments (3a, 3b, and 6b) that aimed at testing whether the moral valence effect found in the previous experiments is modulated by self-referential processes. These three experiments included data from 108 participants.

See Figure 2.

### Part 3: Implicit binding between valence and identity

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

For the task relevant part, we found self-related conditions were performed better than other-related conditions, on both $d$ prime and reaction times.

Most importantly, we found evidence, albeit weak, that task-irrelevant moral valence also played an role. The $d$ prime is greater when shapes were associated with good self condition than with neutral self (BF = 4.4) or bad self (3.1), but shapes associated with bad self and neutral self didn't show differences. In contrast the $d$ prime was smaller when shapes were associated with good other than with neutral other or bad other. See Figure 3.

In this task, we found shapes associated with good person conditions were performed better than other-related conditions, on both $d$ prime and reaction times.

Most importantly, we found evidence, that task-irrelevant self-relevance also played an role. For shapes associated with good person, the $d$ prime was greater when shapes had an "self" inside as task-irrelevant stimuli than with "other" inside (mean_diff = 0.14, 95% credible intervals [-0.02, 0.31], BF = 12.07, p = 0.92), but this effect did not happen when the target shape where associated with "neutral" (mean_diff = 0.04, 95% CI [-.11, .18]) or "bad" person (mean_diff = -.05, 95% CI[-.18, .09]). The same trend appear for the RT

547 data. For shapes associated with good person, an "self" inside will reduce the RTs as

548 compared with when a "other" inside the shape (mean_diff = -55 ms, 95%CI[-75, -35], p <

549 0.0001), but this effect did not occure when the shapes were associated neutral (mean_dfiff

550 = 10, 95% CI [1, 20]) or bad (mean_diff = 5, 95%CI [-16, 27]) person. See Figure 4.

**Self-reported personal distance**

552     See Figure **??**.

**Correlation analyses**

554     The reliability of questionnaires can be found in (Liu et al., 2020). We calculated the

555 correlation between the data from behavioral task and the questionnaire data. First, we

556 calculated the score for each scale based on their structure and factor loading, instead of

557 sum score (McNeish & Wolf, 2020). Then, we used SEM to estimate the correlation

558 because it can include measurement model and statistical model in a unified framework.

559     To make sure that what we found were not false positive, we used two method to

560 ensure the robustness of our analysis. first, we split the data into two half: the data with

561 self and without, then, we used the conditional random forest to find the robust correlation

562 in the exploratory data (with self reference) that can be replicated in the confirmatory data

563 (without the self reference). The robust correlation were then analyzed using SEM

564     Instead of use the exploratory correlation analysis, we used a more principled way to

565 explore the correlation between parameter of HDDM ($v$, $t$, and $a$) and scale scores and

566 person distance.

567     We didn't find the correlation between scale scores and the parameters of HDDM,

568 but found weak correlation between personal distance and the parameter estimated from

569 Good and neutral conditions.

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

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

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

## Discussion

## References

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

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

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

Bagnis, A., Celeghin, A., Mosso, C. O., & Tamietto, M. (2019). Toward an

594  integrative science of social vision in intergroup bias. *Neuroscience &*

595  *Biobehavioral Reviews*, *102*, 318–326.

596  https://doi.org/https://doi.org/10.1016/j.neubiorev.2019.04.020

597  Baumeister, R. F., Bratslavsky, E., Finkenauer, C., & Vohs, K. D. (2001). Bad is

598  stronger than good. *Review of General Psychology*, *5*(4), 323–370.

599  https://doi.org/10.1037/1089-2680.5.4.323

600  Brainard, D. H. (1997). The psychophysics toolbox [Journal Article]. *Spatial Vision*,

601  *10*(4), 433–436.

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

603  stan [Journal Article]. *Journal of Statistical Software; Vol 1, Issue 1 (2017)*.

604  Retrieved from

605  https://www.jstatsoft.org/v080/i01%20http://dx.doi.org/10.18637/jss.v080.i01

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

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

608  2100. https://doi.org/10.1038/s41467-020-15602-4

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

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

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

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

613  Dang, J., King, K. M., & Inzlicht, M. (2020). Why are self-report and behavioral

614  measures weakly correlated? *Trends in Cognitive Sciences*, *24*(4), 267–269.

615  https://doi.org/10.1016/j.tics.2020.01.007

616  Eiserbeck, A., & Abdel Rahman, R. (2020). Visual consciousness of faces in the

617  attentional blink: Knowledge-based effects of trustworthiness dominate over

618  appearance-based impressions. *Consciousness and Cognition*, *83*, 102977.

619  https://doi.org/10.1016/j.concog.2020.102977

Ellemers, N., Toorn, J. van der, Paunov, Y., & Leeuwen, T. van. (2019). The psychology of morality: A review and analysis of empirical studies published from 1940 through 2017. *Personality and Social Psychology Review, 23*(4), 332–366. https://doi.org/10.1177/1088868318811759

Enock, F. E., Hewstone, M. R. C., Lockwood, P. L., & Sui, J. (2020). Overlap in processing advantages for minimal ingroups and the self. *Scientific Reports, 10*(1), 18933. https://doi.org/10.1038/s41598-020-76001-9

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

Freeman, J. B., & Ambady, N. (2011). A dynamic interactive theory of person construal. *Psychological Review, 118*(2), 247–279. https://doi.org/10.1037/a0022327

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 [Journal Article]. *Social and Personality Psychology Compass, 10*(10), 535–549. https://doi.org/10.1111/spc3.12267

Goodwin, G. P. (2015). Moral character in person perception. *Current Directions in Psychological Science, 24*(1), 38–44. https://doi.org/10.1177/0963721414550709

Goodwin, G. P., Piazza, J., & Rozin, P. (2014). Moral character predominates in person perception and evaluation. *Journal of Personality and Social Psychology, 106*(1), 148–168. https://doi.org/10.1037/a0034726

Henrich, J., Heine, S. J., & Norenzayan, A. (2010). The weirdest people in the world? *Behavioral and Brain Sciences*, *33*(2), 61–83. https://doi.org/10.1017/S0140525X0999152X

Hofmann, W., Wisneski, D. C., Brandt, M. J., & Skitka, L. J. (2014). Morality in everyday life. *Science*, *345*(6202), 1340–1343. https://doi.org/10.1126/science.1251560

Hu, C.-P., Lan, Y., Macrae, C. N., & Sui, J. (2020). Good me bad me: Does valence influence self-prioritization during perceptual decision-making? [Journal Article]. *Collabra: Psychology*, *6*(1), 20. https://doi.org/10.1525/collabra.301

Jones, B. C., DeBruine, L. M., Flake, J. K., Liuzza, M. T., Antfolk, J., Arinze, N. C., … Coles, N. A. (2021). To which world regions does the valence–dominance model of social perception apply? *Nature Human Behaviour*, 1–9. https://doi.org/10.1038/s41562-020-01007-2

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

Kirby, K. N., & Gerlanc, D. (2013). BootES: An r package for bootstrap confidence intervals on effect sizes. *Behavior Research Methods*, *45*(4), 905–927. https://doi.org/10.3758/s13428-013-0330-5

Klein, N., & Epley, N. (2016). Maybe holier, but definitely less evil, than you: Bounded self-righteousness in social judgment. *Journal of Personality and Social Psychology*, *110*(5), 660–674. https://doi.org/10.1037/pspa0000050

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 Article]. *Journal of Open Psychology Data*, *8*(1), 1. https://doi.org/10.5334/jopd.49/

Matzke, D., & Wagenmakers, E.-J. (2009). Psychological interpretation of the ex-gaussian and shifted wald parameters: A diffusion model analysis. *Psychonomic Bulletin & Review, 16*(5), 798–817. https://doi.org/10.3758/PBR.16.5.798

McHugh, C., McGann, M., Igou, E. R., & Kinsella, E. (2019). *Moral judgment as categorization (MJAC)*. PsyArXiv. https://doi.org/10.31234/osf.io/72dzp

McNeish, D., & Wolf, M. G. (2020). Thinking twice about sum scores. *Behavior Research Methods*. https://doi.org/10.3758/s13428-020-01398-0

Monin, B., & Jordan, A. H. (2009). The dynamic moral self: A social psychological perspective. In *Personality, identity, and character: Explorations in moral psychology* (pp. 341–354). New York, NY, US: Cambridge University Press. https://doi.org/10.1017/CBO9780511627125.016

Pelli, D. G. (1997). The VideoToolbox software for visual psychophysics: Transforming numbers into movies [Journal Article]. *Spatial Vision, 10*(4), 437–442.

Reicher, S., & Hopkins, N. (2016). Perception, action, and the social dynamics of the variable self. *Psychological Inquiry, 27*(4), 341–347. https://doi.org/10.1080/1047840X.2016.1217584

Rouder, J. N., & Lu, J. (2005). An introduction to bayesian hierarchical models with an application in the theory of signal detection [Journal Article]. *Psychonomic Bulletin & Review, 12*(4), 573–604. https://doi.org/10.3758/bf03196750

Rousselet, G. A., & Wilcox, R. R. (2019). Reaction times and other skewed distributions: Problems with the mean and the median [Preprint]. *Meta-Psychology*. https://doi.org/10.1101/383935

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., & De Brigard, F. (2019). Moral memories and the belief in the good self. *Current Directions in Psychological Science*, *28*(4), 387–391. https://doi.org/10.1177/0963721419847990

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

Stolier, R. M., & Freeman, J. B. (2016). Functional and temporal considerations for top-down influences in social perception. *Psychological Inquiry*, *27*(4), 352–357. https://doi.org/10.1080/1047840X.2016.1216034

Strohminger, N., Knobe, J., & Newman, G. (2017). The true self: A psychological concept distinct from the self: *Perspectives on Psychological Science*. https://doi.org/10.1177/1745691616689495

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

Tappin, B. M., & McKay, R. T. (2017). The illusion of moral superiority. *Social Psychological and Personality Science*, *8*(6), 623–631. https://doi.org/10.1177/1948550616673878

Turner, J. C., Hogg, M. A., Oakes, P. J., Reicher, S. D., & Wetherell, M. S. (1987). *Rediscovering the social group: A self-categorization theory.* Cambridge, MA, US: Basil Blackwell.

Turner, J. C., Oakes, P. J., Haslam, S. A., & McGarty, C. (1994). Self and collective: Cognition and social context. *Personality and Social Psychology Bulletin*, *20*(5), 454–463. https://doi.org/10.1177/0146167294205002

Uhlmann, E. L., Pizarro, D. A., & Diermeier, D. (2015). A person-centered approach to moral judgment: *Perspectives on Psychological Science.* https://doi.org/10.1177/1745691614556679

Verosky, S. C., & Todorov, A. (2010). Differential neural responses to faces physically similar to the self as a function of their valence. *NeuroImage*, *49*(2), 1690–1698. https://doi.org/10.1016/j.neuroimage.2009.10.017

Waytz, A., Dungan, J., & Young, L. (2013). The whistleblower's dilemma and the fairness–loyalty tradeoff. *Journal of Experimental Social Psychology*, *49*(6), 1027–1033. https://doi.org/10.1016/j.jesp.2013.07.002

Wiecki, T. V., Sofer, I., & Frank, M. J. (2013). HDDM: Hierarchical bayesian estimation of the drift-diffusion model in python. *Frontiers in Neuroinformatics*, *7.* https://doi.org/10.3389/fninf.2013.00014

Willis, J., & Todorov, A. (2006). First impressions: Making up your mind after a 100-ms exposure to a face. *Psychological Science*, *17*(7), 592–598. https://doi.org/10.1111/j.1467-9280.2006.01750.x

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
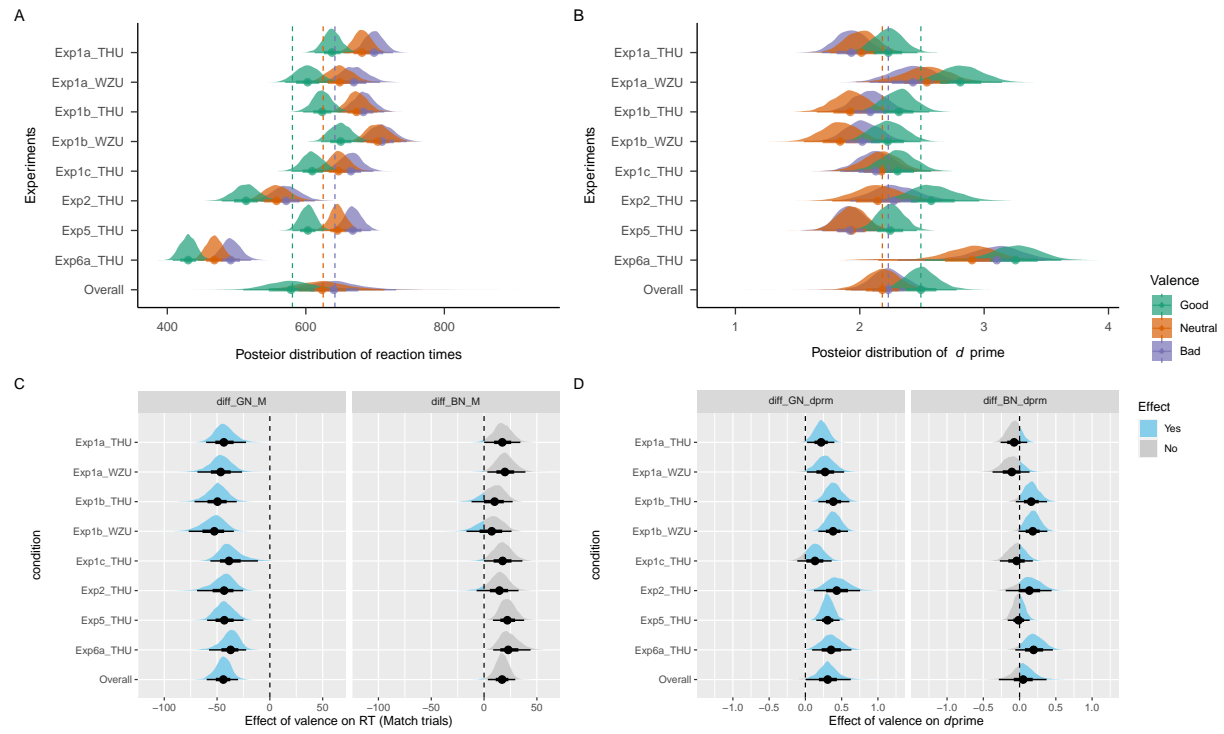
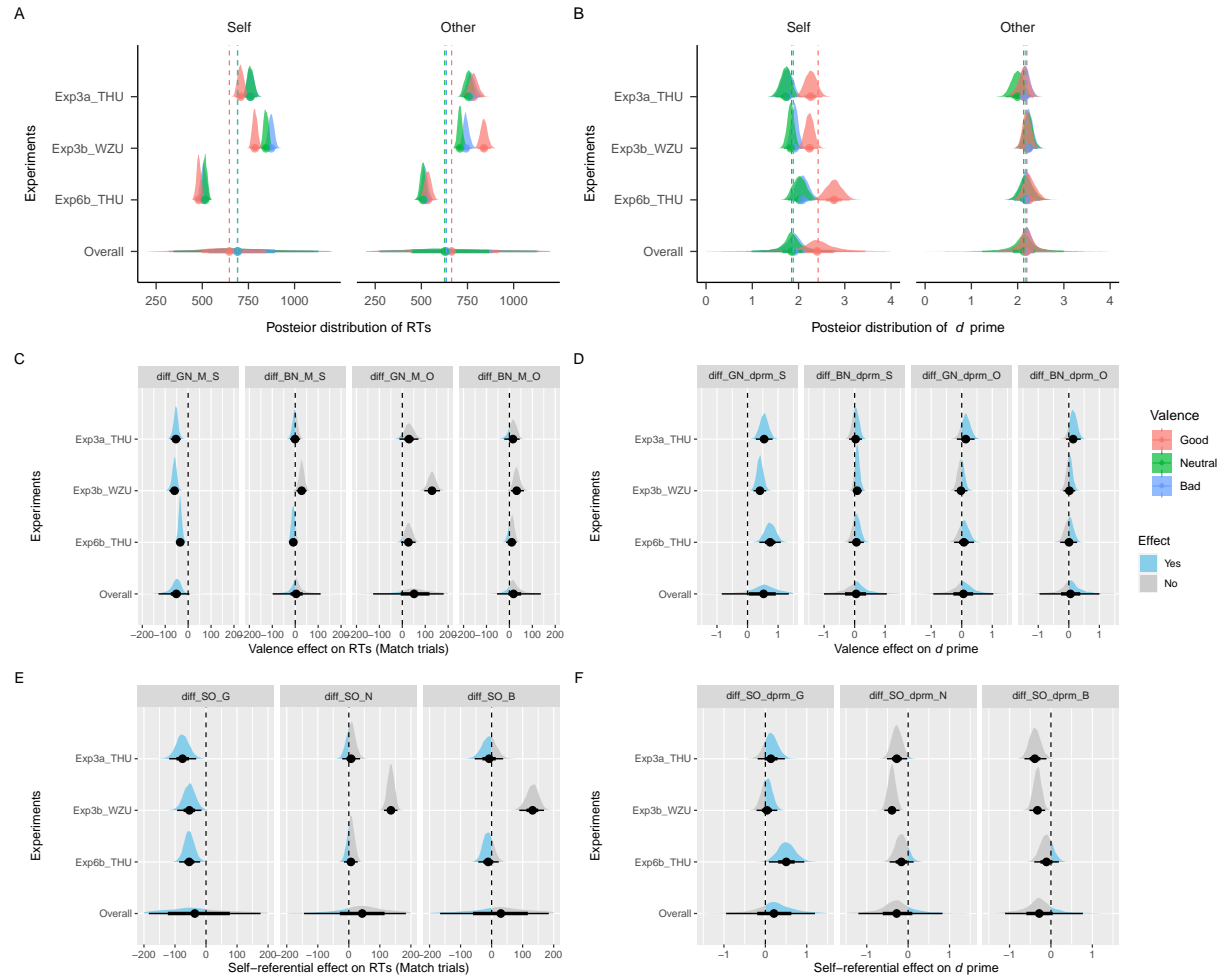*Figure 1.* Effect of moral valence on RT and d'

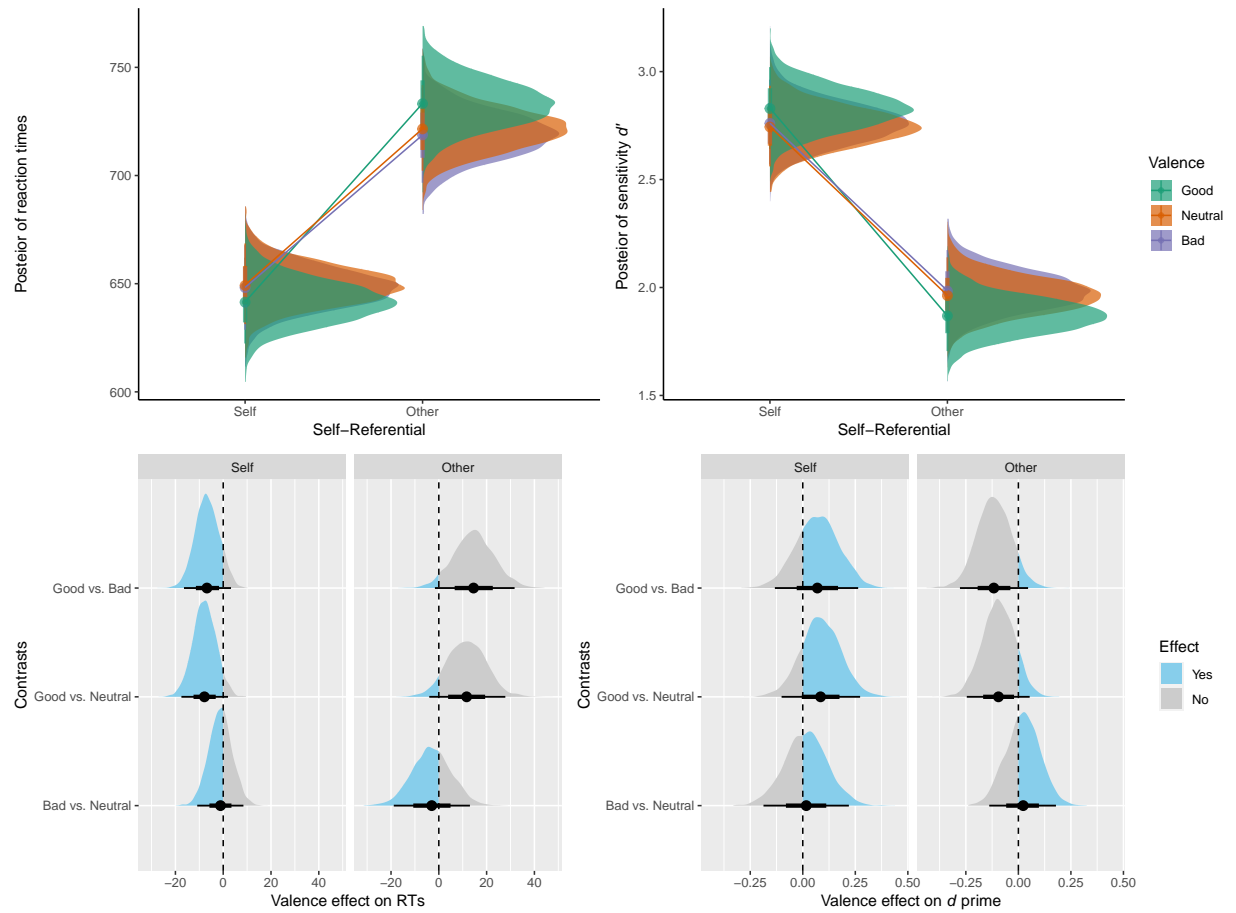*Figure 2.* Interaction between moral valence and self-referential
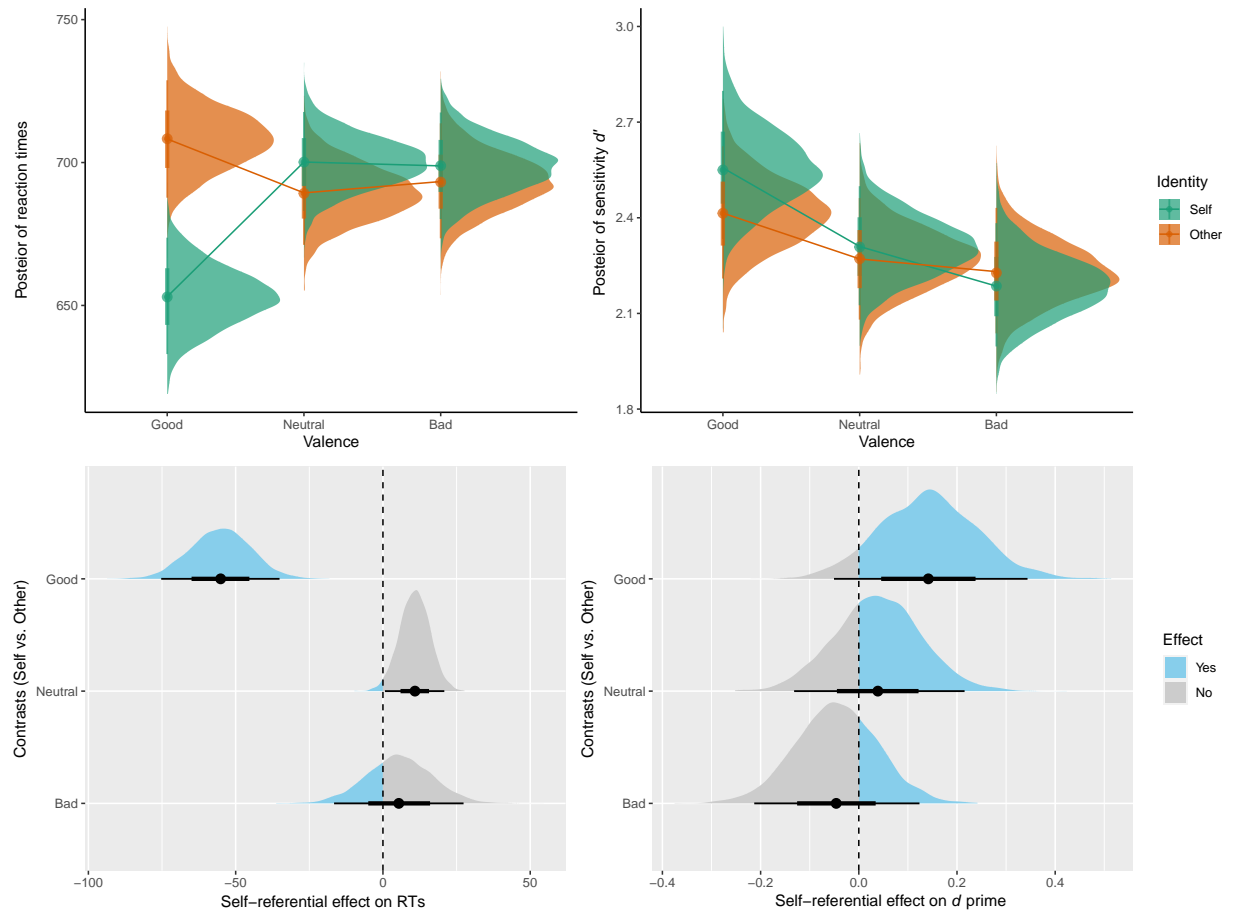
*Figure 3*. exp4a: Results of Bayesian GLM analysis.

*Figure 4.* exp4a: Results of Bayesian GLM analysis.