**Title**

An empirical validation of the Quadruple Process Model of implicit attitudes against alternate, theoretically defensible specifications.

**Permalink**

https://escholarship.org/uc/item/6mk381s0

**Authors**

Calanchini, Jimmy
Sherman, Jeff
Klauer, Karl Christoph
et al.

**Publication Date**

2017-02-15

**DOI**

10.31234/osf.io/7b6nz

Peer reviewed

An Empirical Validation of the Quadruple Process Model of Implicit Attitudes Against

Alternate, Theoretically Defensible Specifications.

Jimmy Calanchini[1*], Jeffrey W. Sherman[1], Karl Christoph Klauer[2], Emilio Ferrer[1]

[1] Institut für Psychologie, Albert Ludwig Universitat, Freiburg, Germany

[2] Department of Psychology, University of California, Davis, USA

*Corresponding author

E-mail: jcalanchini@ucdavis.edu

Last revised: 15 February 2017

**Abstract**

The Quadruple process (Quad) model is a multinomial processing tree that specifies the joint contribution of four qualitatively distinct cognitive processes to responses on implicit measures. The way in which these processes interact to drive responses was initially specified according to theory, and the construct validity of this specification of the model has been demonstrated across a wide variety of studies. However, there are other theoretically-defensible ways in which these processes might interact. The purpose of the present research was to compare the standard version of the Quad model against alternate specifications in order to determine which model best fits data from the Implicit Association Test. Three different versions of the Quad model were applied to very large samples of real participants' data across three content domains: racial attitudes, sexual orientation attitudes, and gender stereotypes. The standard model provided best fit for racial attitudes and gender stereotype data. However, other versions of the model provided equivalent fit to sexual orientation attitudes data. Taken together, these analyses indicate that the standard version of the Quad model provides best fit to data from the Implicit Association Test in general, but that alternate specifications may be appropriate for some content domains and participant populations.

KEYWORDS: implicit attitudes; prejudice; stereotyping; multinomial model; Quad model

## Introduction

Implicit attitude measures were created to overcome problems associated with self-report (or explicit) attitude measures that had troubled researchers for decades. Explicit measures, which directly ask respondents to report their attitudes, are susceptible to deliberate response strategies that may arise from social desirability or self-presentational concerns. Explicit measures also are unable to capture mental content that is inaccessible through introspection [1]. In contrast, implicit measures were designed to minimize these "unwilling and unable" problems by assessing attitudes and beliefs without directly requesting that respondents report those attitudes and beliefs. Implicit measures indirectly assess attitudes in a variety of ways, such as structuring the task in a manner that conceals what is being measured (e.g., evaluative priming [2]; semantic priming [3]) or by making responses difficult to control (e.g., IAT [4]; GNAT [5]).

These features of implicit measures initially led to the widely-held belief that responses on these measures reflect only the respondent's underlying and automatically-activated mental associations with the attitude object (e.g.,[6, 4]). However, subsequent research employing mathematical modeling techniques revealed that responses on implicit measures are not process pure but, instead, reflect the joint contribution of multiple processes. The focus of the present research is on the Quadruple Process model (Quad model: [7, 8]), though a number of other such models have also identified the influence of multiple processes on implicit measures (e.g., [9, 10, 11, 12, 13, 14). The Quad model made an important contribution to the field of social psychology by specifying the influence of four qualitatively distinct processes to responses on implicit measures: the activation of biased associations, an accuracy-oriented process, an inhibitory process, and a bias that drives responses in the absence of other guides.  Conrey and colleagues [7] demonstrated the stochastic and construct validity of the Quad model, and it has

since provided excellent fit to a wide variety of empirical data. Table 1 lists the implicit

measures and content domains to which the Quad model has been successfully applied to date.

However, just because the Quad model fits data well does not necessarily mean that it provides

the best possible fit.  The precise manner in which these four processes are specified to interact

was initially articulated based on theory, but there are other theoretically-defensible ways in

which these processes might interact, and an alternate specification might provide superior fit.

Thus, the purpose of the present research is to compare the standard version of the Quad model

against alternate, theoretically-defensible specifications in order to determine which model best

fits data from the IAT, the most commonly used implicit measure and the measure to which the

Quad model has been most frequently applied.

**Table 1.**

| Citation | Measure | Domain | Population |
|---|---|---|---|
| Jin et al., (2016) | IAT | male-female evaluative | undergraduate (China) |
| Scroggins et al.,  (2015) | IAT | Black-White / ingroup-outgroup evaluative | undergraduate (USA) |
| Huntsinger et al., (2015) | IAT | Black-White evaluative | undergradaute (USA) / mTurk |
| Burke (2015) | WIT | Black-White / gun-tool | simulated |
| Ramos et al., (2015) | GNAT | male-female career-family | undergraduate (Portugal) |
| Zestcott et al., (2015) | IAT | tattoo evaluative | undergraduate (USA) |
| Ruiz et al., (2015) | IAT | old-young evaluative | medical students (USA) |
| Calanchini et al., (2014) | | | |
| Study 1a | IAT | Black-White evaluative, Asian-White evaluative | undergraduate (USA) |
| Study 1b | IAT | Black-White evaluative, flower-insect evaluative | undergraduate (USA) |

| | | | |
|---|---|---|---|
| Study 1c | IAT | Black-White physical-mental, flower-insect evaluative | undergraduate (USA) |
| Study 2a | IAT | Black-White evaluative, skin tone evaluative | Project Implicit |
| Study 2b | IAT | gay-straight evaluative, disability evaluative | Project Implicit |
| Study 2c | IAT | old-young evaluative, male-female career-family | Project Implicit |
| Clerkin et al., (2014) | IAT | me-not me / calm-panicked | clinical (USA) |
| Gonsalkorale et al., (2013) | IAT | old-young evaluative | Project Implicit |
| Calanchini et al., (2013) | IAT | Black-White evaluative | undergraduate (USA) |
| Soderberg & Sherman (2013) | IAT | Black-White evaluative | undergraduate (USA) |
| O'Connor et al., (2012) | single-category IAT | alcohol / cigarette evaluative | children (USA) |
| Allen & Sherman (2011) | IAT | Black-White evaluative | undergraduate (USA) |
| Gonsalkorale et al., (2011) | | | |
| Study 1 | WIT | Black-White / gun-tool | undergraduate (USA) |
| Study 2 | IAT | Black-White evaluative | undergraduate (USA) |
| Allen et al., (2010) | priming | Black-White evaluative | undergraduate (USA) |
| Gonsalkorale et al., (2010) | IAT | Black-White evaluative | Project Implicit / undergraduate (USA) |
| Gonsalkorale et al., (2009a) | IAT | Black-White evaluative | Project Implicit |
| Gonsalkorale et al., (2009b) | GNAT | Muslim-White evaluative | Undergraduate (Australia) |
| Beer et al., (2008) | IAT | Black-White evaluative | undergraduate (USA) |
| Bishara & Payne (2009) | WIT | Black-White gun-tool | undergraduate (USA) |
| Conrey et al., (2005) | | | |
| Study 1 | IAT | flower-insect | undergraduate (USA) |

| | | evaluative | |
|---|---|---|---|
| Study 2 | IAT | Black-White | undergraduate (USA) |
| | | evaluative | |
| Study 3 | IAT | flower-insect | undergraduate (USA) |
| | | evaluative | |
| Study 4 | IAT | Black-White | undergraduate (USA) |
| | | evaluative | |
| Study 5 | IAT | Black-White | undergraduate (USA) |
| | | evaluative | |

*Note*: IAT: Implicit Association Test. WIT: Weapons Identification Task. GNAT: Go/No-Go Association Task.

## The Quad Model

Though the Quad model has been applied most extensively to the IAT [4], it can be applied to most tasks that are based on the logic of response compatibility (cf. [15, 16]). For example, an IAT designed to measure attitudes toward Black relative to White people might present pictures of Black and White people and pleasant and unpleasant words. On some trials, participants press one button in response to pictures of White people and pleasant words and another button in response to pictures of Black people and unpleasant words. On other trials, the response labels are changed such that participants press one button in response to pictures of White people and unpleasant words and another button in response to pictures of Black people and pleasant words. The ease with which participants can respond to one set of pairings relative to the other set of pairings (measured in terms of response latency or accuracy) has generally been interpreted as reflecting relative differences in how strongly the concepts are associated in memory. That is, if a participant responds more quickly or accurately when White and pleasant stimuli share a response key than when Black and pleasant stimuli share a response key, then she is assumed to associate pleasant concepts more strongly with White people than with Black people.

Building upon this interpretation of responses on implicit measures, the Quad model specifies that associations stored in memory are just one of several processes that contribute to implicit task performance. Specifically, the Quad model posits the influence of Activation of Associations (AC), Detection of correct responses (D), Overcoming Bias (OB), and Guessing (G). The AC parameter refers to the degree to which biased associations are activated when responding to a stimulus. All else equal, the stronger the associations, the more likely they are to be activated and to drive behavior in an association-consistent direction. The D parameter reflects the likelihood that the participant can accurately discriminate between correct and incorrect responses. Sometimes, the activated associations conflict with the detected correct response. For example, when the categories "Black" and "pleasant" share a response key, activated biased racial associations (e.g., between Black and unpleasant) might conflict with the correct response (i.e., to press the same button for Black and pleasant stimuli). In such cases, the Quad model proposes that an OB process resolves the conflict. As such, the OB parameter refers to an inhibitory process that prevents activated associations from influencing behavior when they conflict with detected correct responses. Finally, the G parameter reflects biases that drive responses when the participant has no associations that direct behavior and is unable to detect the correct response.

The Quad model has been instantiated as a multinomial processing tree [17, 18]. A portion of the Quad model depicted as a processing tree is presented in Figure 1. In the tree, each path represents a likelihood. Processing parameters with lines leading to them are conditional on all preceding parameters. For instance, OB is conditional on both AC and D. The conditional relationships described by the model form a system of equations that predicts the numbers of correct and incorrect responses to different stimulus pairings on the implicit measure. Note that

these conditional relationships do not imply a serial or temporal order in the onset and

conclusion of the different processes. Rather, these relationships are mathematical descriptions

of the manner in which the parameters interact to produce behavior. Thus, the activation of

associations (AC), attempts to detect a correct response (D), and attempts to overcome

associations (OB) may occur simultaneously. However, in determining a response on an

incompatible trial, the status of OB determines whether AC or D drives responses when they are
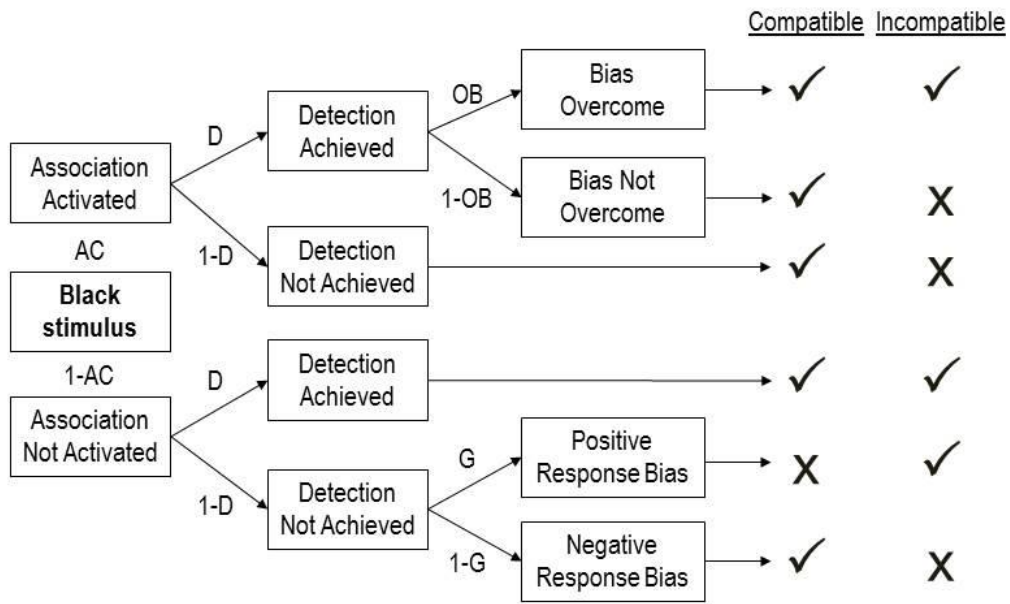
in conflict.



*Fig 1.*

A portion of the Quadruple Process Model (Quad Model). Each square represents a parameter

and each path represents a likelihood. All parameters are conditional upon all preceding paths.

The table on the right side of the figure depicts correct (✓) and incorrect (✗) responses as a function of process pattern.

Here is an example of how the conditional relationships described by the processing tree in Figure 1 spell out the equations that comprise the Quad model. In the standard version of the model there are three ways in which an incorrect response can be returned on a trial in which "Black" and "pleasant" share a response key. The first is that biased associations are activated (AC), detection succeeds (D), and OB fails $(1 - OB)$, the likelihood of which can be represented by the equation $AC \times D \times (1 - OB)$. The second is that the biased associations are activated (AC) and detection fails $(1 - D)$, with probability given by the equation $AC \times (1 - D)$. The third is the likelihood that biased associations are not activated $(1 - AC)$, detection fails $(1 - D)$, and a bias toward guessing "unpleasant" $(1 - G)$ produces an incorrect response, with likelihood represented by the equation $(1 - AC) \times (1 - D) \times (1 - G)$. As such, the overall likelihood of producing an incorrect response on this trial is represented as the sum of these three conditional probabilities: $[AC \times D \times (1 - OB)]+[AC \times (1 - D)]+[(1 - AC) \times (1 - D) \times (1 - G)]$. The respective equations for each item category (e.g., White, Black, pleasant words, and unpleasant words) in each trial type (e.g., compatible; incompatible) are then used to predict the observed proportions of errors in a given data set. The model's predictions are compared to the actual data to determine the model's ability to account for the data. A model fit estimate (e.g., chi-square) is computed for the difference between the predicted and observed responses. To best approximate the model to the data, the parameter values are changed until they produce a minimum possible

value of the model fit statistic. The final parameter values that result from this process are interpreted as relative levels of the processes.

**Model 1: The standard version of the Quad model.** Several important assumptions are built into the standard version of the Quad model. The first assumption is that certain stimulus pairings on implicit measures are easier to respond to (e.g., White-pleasant) than other stimulus pairings (e.g., White-unpleasant). Trials that consist of pairings that are easier for most participants to respond to (measured in terms of either response latency or accuracy) are normatively referred to as *compatible* trials, whereas trials that are more difficult for participants to respond to are normatively referred to as *incompatible* trials. Scores of empirical data support this assumption of stimulus compatibility: for example, most people who complete IATs on the Project Implicit website (https://implicit.harvard.edu) can more quickly and accurately respond to White-pleasant pairings than Black-pleasant pairings [19].

The Quad model accounts for the fact that some stimulus pairings are easier to respond to than others by assuming that different combinations of processes influence responses on each trial type. On a compatible IAT trial (e.g., White and pleasant share a response key), the Quad model specifies that both accuracy-oriented detection (D) and activated White-pleasant associations (AC) produce the same response tendency: to press the button labeled White/pleasant. Thus, there is no need for the inhibitory OB process to intervene on a compatible trial because there are no conflicting response tendencies that need to be resolved. In contrast, on an incompatible IAT trial (e.g., White and unpleasant share a response key), accuracy-oriented detection (D) produces one response tendency (i.e., to press the button labeled White/unpleasant) but activated White-pleasant associations (AC) produce a different response tendency (i.e., to press the button labeled Black/pleasant). In this case, the OB process must intervene on an

incompatible trial to resolve the discrepancy between accuracy-oriented detection and activated associations. If OB succeeds, then the participant will make the correct response but, if OB fails, the participant will make the incorrect response. Thus, the Quad model assumes that the OB process is only active on incompatible trials.

The second assumption built into the Quad model is that OB influences responses to *target* stimuli (e.g., pictures of Black and White people) but not *attribute* stimuli (e.g., pleasant and unpleasant words) on incompatible trials. In other words, biased associations only need to be inhibited when responding to stimuli representing the target groups. This assumption is based, in part, on the idea that associative links are not always bidirectional: for example, a picture of a White person may activate pleasant concepts, but a pleasant word may not necessarily activate White concepts. However, in the original specification of the Quad model [7], this was not the case: two separate OB parameters were estimated on incompatible trials, one that influenced responses to target stimuli and another that influenced responses to attribute stimuli. The current specification of the Quad model is based on the expectation that the OB process is limited to social targets, and is not activated by attribute stimuli. Because OB does not intervene on responses to attribute stimuli on incompatible trials, these trials are essentially a conflict between activated associations (AC) and accuracy-oriented detection (D). An important assumption of the model is that, on such trials, activated associations will drive (incorrect) responses on incompatible attribute trials and that accuracy-oriented Detection can only drive a correct response on such trials when biased associations are not activated.

The full set of equations that comprise the standard version of the Quad model can be found in the Appendix. The assumptions that OB only influences responses to target stimuli on incompatible trials and that detection (D) can only drive responses to attribute stimuli on

incompatible trials when biased associations are not activated have not yet been tested empirically against their logical alternatives. It is possible to specify a model in which OB influences responses to both incompatible target and attribute stimuli. Similarly, it is possible to specify a model in which biased associations (AC) can only drive responses on incompatible trials when detection (D) fails (in contrast to the current model, in which D can only drive responses on incompatible attribute trials when there is no AC).  In the present research, we specified these alternative models and empirically tested them against the standard version of the Quad model (hereafter referred to as Model 1).

**Model 2: Control-dominant.** In the standard version of the Quad model (Model 1), the relatively more automatic process (i.e., activated associations, AC) dominates responses on incompatible attribute trials. That is, the relatively more controlled process (detection, D) can only drive responses when AC is not activated. This mirrors the relationship between automatic and controlled processing proposed by Lindsay and Jacoby [20] to account for performance on the Stroop task. Alternately, it is also theoretically possible that the relatively more controlled process dominates responses on incompatible trials. Thus, Model 2 specifies that accuracy-oriented detection (D) drives responses on incompatible attribute trials and that activated associations (AC) can only drive responses when detection fails. This mirrors the relationship between automatic and controlled processing proposed by Jacoby [21] to account for source memory performance (see also [12]).The processing tree for Model 2 can be found in the Appendix.

**Model 3: OB drives responses on all incompatible trials.** In the standard version of the model (Model 1), OB only influences responses to *target* stimuli (e.g., Black and White faces) on incompatible trials. Alternately, it is possible that OB also influences responses to *attribute*

stimuli (e.g., pleasant and unpleasant words) on incompatible trials. This version of the model (referred to as Model 3) can be found in the Appendix, and specifies that, when activated associations and accuracy-oriented Detection would produce different responses (i.e., on incompatible trials), OB intervenes to resolve this discrepancy for both target and attribute stimuli.

**Overview of Analyses**

The purpose of the present analyses is to empirically compare these three different specifications of the Quad model. Model 1 is the standard version of the model, in which OB only influences responses to incompatible target trials and detection (D) can only drive responses on incompatible attribute trials when biased associations (AC) are not activated.  Model 2 reverses the relation between AC and D on incompatible attribute trials: biased associations (AC) can only drive responses on incompatible attribute trials when detection (D) fails.  In Model 3, OB resolves any discrepancies between AC and D on both incompatible target and attribute trials and, as such, obviates the issue of whether D or AC dominates.

We applied the three different versions of the model to real participant's IAT data from three different content domains. The first IAT assessed relations between Black and White people and positive and negative concepts (i.e., racial attitudes). This test is arguably the most commonly-used IAT, so it is ecologically valid to assess the different versions of the Quad model on such a ubiquitous test. This also is the content domain to which the Quad model has been most frequently applied (Table 1). The second IAT assessed relations between gay and straight people and positive and negative concepts (i.e., sexuality attitudes). This is also a widely-used test, though not nearly as commonly used as the Black/White IAT.  The third IAT assessed relations between males and females and concepts related to careers and families (i.e.,

gender stereotypes). Whereas the racial and sexuality IATs assessed relations between social

groups and evaluations (i.e., attitudes), the gender IAT assessed relations between social groups

and attributes (i.e., stereotypes), and was included to ensure that these analyses were not limited

to the evaluative domain.

One final purpose of the present analyses was to assess the robustness of the model by

fitting it to data from participants with non-normative response tendencies. The IAT *d* score [22]

is a latency-based index of the ease with which participants can respond to one set of stimulus

pairings (e.g., White-pleasant) relative to another (e.g., Black-pleasant). In addition to applying

the model to randomly-selected data, we also applied it to data from participants who scored

approximately +1SD and -1SD from the mean on the IAT *d* score in order to observe model fit

for participants for whom certain pairings were especially facilitated or especially hindered.

**Participants, Materials, and Procedure**

All data were collected from visitors to the Project Implicit website

([https://implicit.harvard.edu](https://implicit.harvard.edu)) between 2006 to 2010 who chose to complete one of up to 15

possible IATs [23]. The structure of the IATs was based on that described by Nosek, Banaji, and

Greenwald [24]. Participants first completed a 20-trial practice block in which they categorized

target stimuli (e.g., images of Black and White people) followed by a 20-trial practice block in

which they categorized attribute stimuli (e.g., pleasant and unpleasant words). In the third (20

trial) and fourth (40 trial) critical blocks, participants simultaneously categorized both attribute

and target stimuli. The fifth 40-trial block consisted of categorizing only target stimuli, but with

the response keys reversed. The sixth and seventh critical blocks were identical to the third and

fourth blocks, but the response keys reflected the switched target pairing of the fifth block. If a

participant made an error in categorization during any of the response trials, a red "X" appeared

below the stimulus and remained there until the participant corrected the error.

The racial attitudes IAT consisted of pleasant and unpleasant words (e.g., pleasure,

failure) and images of Black and White males. The attribute category labels were

Pleasant/Unpleasant and the target category labels were Black/White. The full sample of racial

attitudes IAT data consisted of 457,379 participants, 56.7% female, $M_{age}$=26.93, $SD_{age}$=11.36.

The sexuality attitudes IAT consisted of pleasant and unpleasant words (e.g., pleasure,

failure) and words and images representing gay and straight people (e.g., "homosexual,"

"heterosexual," a picture of two brides, a picture of a bride and a groom). The attribute category

labels were Pleasant/Unpleasant and the target category labels were Straight People/Gay People.

The full sample of sexuality attitudes IAT data consisted of 11,896 participants, 73.5% female,

$M_{age}$=24.62, $SD_{age}$=10.33.

The gender stereotypes IAT consisted of male and female names (e.g., John, Michelle)

and words representing careers and families (e.g., office, home). The attribute category labels

were Family/Career and the target category labels were Male/Female. The full sample of gender

stereotypes IAT data consisted of 10,990 participants, 77.6% female, $M_{age}$=24.96, $SD_{age}$=10.61.

**Results & Discussion**

For all three IATs, ten samples of 500 participants each was drawn (with replacement)

from the full sample of participants. Additionally, ten samples of 500 participants each was

drawn (with replacement) from participants who scored approximately 1SD above the mean on

the *d* score, as well as from participants who scored approximately 1SD below the mean on the *d*

score. Given that we had such large datasets (all Ns>10,000), we treated each dataset as a

"population" and drew samples from it in order to simulate sampling variability, providing a

measure of ecological validity to these analyses. Table 2 presents *d* scores and other descriptive

statistics averaged across each set of ten 500-person samples.

**Table 2.**

*Descriptive Statistics*

|  | IAT *d* (SD) | Min. | Max. | Skew (SE) | Kurtosis (SE) | Com. Errors | Inc. Errors | Total Errors |
|---|---|---|---|---|---|---|---|---|
| Black/White |  |  |  |  |  |  |  |  |
| full sample | .35 (.42) | -1.84 | 1.49 | -.47 (.04) | .27 (.07) | 5.99% | 9.64% | 7.81% |
| +1SD | .77 (.02) | 0.72 | 0.81 | .07 (.04) | -1.18 (.07) | 4.60% | 11.26% | 7.94% |
| -1SD | -.09 (.02) | -0.13 | -0.05 | -.09 (.04) | -1.17 (.07) | 7.29% | 7.76% | 7.52% |
| gay/straight |  |  |  |  |  |  |  |  |
| full sample | .31 (.45) | -1.71 | 1.68 | -.42 (.04) | .14 (.07) | 9.79% | 13.69% | 11.74% |
| +1SD | .76 (.03) | 0.72 | 0.81 | .04 (.04) | -1.19 (.07) | 7.58% | 14.60% | 11.09% |
| -1SD | -.14 (.05) | -0.23 | -0.05 | -.12 (.04) | -1.12 (.07) | 11.27% | 11.61% | 11.44% |
| gender/career |  |  |  |  |  |  |  |  |
| full sample | .38 (.35) | -1.37 | 1.34 | -.40 (.04) | .31 (.07) | 6.47% | 9.56% | 8.02% |
| +1SD | .72 (.04) | 0.65 | 0.79 | .17 (.04) | -1.17 (.07) | 5.60% | 10.67% | 8.13% |
| -1SD | .03 (.04) | -0.05 | 0.09 | -.12 (.04) | -1.11 (.07) | 7.56% | 8.52% | 8.04% |

*Note:* Com. Errors = percent of errors on compatible trials. Inc. Errors = percent of errors on incompatible trials.

Each of the three versions of the Quad model was then applied to each sample. Parameter

estimates of AC, D, OB, and G were calculated for each IAT. The G parameter was coded so that

higher scores represented a bias toward guessing with the "pleasant" key on the racial and

sexuality attitudes IATs, and a bias towards guessing with the "career" key on the gender

stereotypes IAT. For each IAT, two AC parameters were estimated, representing the extent to

which associations with the (normatively) higher status group (AC1) and associations with the

(normatively) lower status group (AC2) were activated. For the racial attitudes IAT, AC1

represents White-pleasant associations and AC2 represents Black-unpleasant associations. For

the sexuality attitudes IAT, AC1 represents straight-pleasant associations and AC2 represents

gay-unpleasant associations. For the gender stereotypes IAT, AC1 represents male-career

associations and AC2 represents female-family associations.

**Full sample data.** Averaged model fit indices, model selection indices, and parameter

estimates for the ten 500-person samples of data randomly selected from the full sample are

presented in Table 3. Chi-square values are presented as indices of absolute model fit. However,

chi-square tests are dependent on sample size, such that minute deviations from the model can

jeopardize model fit when power is high [25]. Given that each sample is comprised of 500

participants and each participant completed 120 critical IAT trials, each sample consists of

60,000 observations. As such, $w$ values also are presented, which represent the effect size of lack

of model fit between the actual data and the model's predicted data, controlling for sample size

(see [34] for additional guidelines on interpreting $w$). Additionally, because all three models

estimate the same number of parameters, chi-square values can be compared among models as

indices of relative best fit. Minimum description length (MDL), which takes into account

differences between models in terms of flexibility due to the different functional forms of the

equations, is also presented as a corroborating index of best fit [26]. We do not report the model

selection indices AIC and BIC because they control for the number of parameters in a model but

not the form (i.e,. complexity) of the model. Given that our models vary in complexity but not in

number of parameters, MDL therefore provides a superior index of model selection.

**Table 3.**

*Model fit / selection indices and parameter estimates for full sample data.*

|  | $\chi^2$ | $w$ | MDL | AC1 | AC2 | D | G | OB |
|---|---|---|---|---|---|---|---|---|
| **Black/White** | | | | | | | | |
| 1 | 39.08 | 0.03 | 16,168.20 | 0.08 (.01) | 0.04 (.00) | 0.79 (.09) | 0.56 (.06) | 0.93 (.04) |
| 2 | 273.62 | 0.07 | 16,284.91 | 0.32 (e) | 0.13 (e) | 0.84 (.00) | 0.51 (.01) | 1.00 (e) |
| 2* | 289.65 | 0.07 | 16,290.48 | 0.32 (.05) | 0.12 (.02) | 0.76 (.09) | 0.51 (.01) | Fixed |
| 2† | 273.62 | 0.07 | 16,282.46 | 0.33 (.05) | 0.13 (.02) | 0.84 (.00) | 0.51 (.06) | N/A |
| 3 | 263.12 | 0.07 | 16,279.36 | 0.17 (.07) | 0.07 (.03) | 0.86 (.01) | 0.53 (.01) | 0.77 (.20) |
| **Gay/Straight** | | | | | | | | |
| 1 | 73.19 | 0.03 | 21,578.44 | 0.05 (.00) | 0.02 (.00) | 0.80 (.00) | 0.52 (.01) | 0.00 (.11) |
| 2 | 78.71 | 0.04 | 21,580.64 | 0.20 (.02) | 0.06 (.02) | 0.77 (.00) | 0.49 (.01) | 0.86 (.04) |
| 3 | 73.19 | 0.03 | 21,577.58 | 0.05 (.03) | 0.02 (.01) | 0.80 (.01) | 0.52 (.01) | 0.00 (.75) |
| **Gender/Career** | | | | | | | | |
| 1 | 106.31 | 0.04 | 16,514.16 | 0.04 (.00) | 0.08 (.00) | 0.87 (.00) | 0.48 (.01) | 0.96 (.05) |
| 2 | 345.63 | 0.08 | 16,633.26 | 0.10 (e) | 0.28 (e) | 0.84 (.00) | 0.53 (.01) | 1.00 (e) |
| 2* | 357.56 | 0.08 | 16,636.77 | 0.08 (.02) | 0.27 (.02) | 0.84 (.00) | 0.53 (.01) | Fixed |
| 2† | 345.63 | 0.08 | 16,630.81 | 0.10 (.02) | 0.28 (.02) | 0.84 (.00) | 0.53 (.01) | N/A |
| 3 | 311.82 | 0.07 | 16,616.05 | 0.02 (.00) | 0.05 (e) | 0.86 (.00) | 0.51 (.01) | 0.04 (.85) |
| 3‡ | 311.83 | 0.07 | 16,614.13 | 0.02 (.00) | 0.05 (.00) | 0.86 (.00) | 0.51 (.01) | Fixed |

*Note*: AC1 refers to associations between the (normatively) higher status group and the (normatively) more positive attribute, and AC2 refers to associations between the (normatively) lower status group and the (normatively) more negative attribute. For the Black/White IAT, AC1 represents White-pleasant associations and AC2 represents Black-unpleasant associations. For the gay/straight IAT, AC1 represents straight-pleasant associations and AC2 represents gay-unpleasant associations. For the gender/career IAT, AC1 represents male-career associations and AC2 represents female-family associations. *OB fixed @ 0.99 , †OB deleted, ‡OB fixed @ 0.01. (Standard Errors). (e)=unable to estimate standard error.

For the racial attitudes IAT, Model 1 ($\chi^2$=39.08, MDL=16,168.20) provided better fit by wide margins across model selection indices than Model 2 ($\chi^2$=273.62, MDL=16,284.91) and Model 3 ($\chi^2$=263.12, MDL=16,279.36). Several issues arose when Model 2 was fit to these data: the model was unable to estimate standard errors for AC1 (i.e., White-pleasant associations), AC2 (i.e., Black-unpleasant associations), and OB, and estimated the OB parameter at ceiling.

Fixing the OB parameter at 0.99 ($\chi^2$=289.65, MDL=16,290.48), as well as deleting it from the model entirely ($\chi^2$=273.62, MDL=16,282.46), allowed standard errors to be estimated, but did not improve model fit.

For the sexuality attitudes IAT data, Model 1 ($\chi^2$=73.19, MDL=21,578.44) and Model 3($\chi^2$=73.19, MDL=21,577.58) provided nearly identical fit across model selection indices, and both provided better fit than Model 2 ($\chi^2$=78.71, MDL=21,580.64). Interestingly, both Models 1 and 3 estimated identical parameters, including zero estimates for the OB parameter, suggesting that this inhibitory process does not significantly contribute to responses on the gay/straight IAT. Because Models 1 and 3 collapse to the same model when OB is zero, these data are not diagnostic for discriminating between them, but can be taken as evidence in favor of these models relative to Model 2.

For the gender stereotypes IAT data, Model 1 ($\chi^2$=106.31, MDL=16,514.16) provided better fit by wide margins across model selection indices than Model 2 ($\chi^2$=345.63, MDL=16,633.26) and Model 3 ($\chi^2$=311.82, MDL=16,616.05).  Several issues arose when Model 2 was fit to these data: the model was unable to estimate standard errors for AC1 (i.e., male-career associations), AC2 (i.e., female-family associations), and OB, and estimated the OB parameter at ceiling. Fixing the OB parameter at 0.99 ($\chi^2$=357.56, MDL=16,636.77), as well as deleting it from the model entirely ($\chi^2$=345.63, MDL=16,630.81), allowed standard errors to be estimated, but did not improve model fit. Similarly, Model 3 was unable to estimate standard errors for AC2. Fixing the OB parameter to 0.01 ($\chi^2$=311.83, MDL=16,614.13) allowed standard errors to be estimated, but did not improve model fit.

+**1SD data.** Averaged model fit indices, model selection indices, and parameter estimates for the ten 500-person samples of data randomly selected from approximately 1SD

above the mean on the IAT $d$ score for each test are presented in Table 4. A total of 27,019

participants demonstrated implicit bias (as indexed by the $d$ score) approximately 1SD above the

mean ($0.9<Z_{IAT}<1.1$) on the racial attitudes IAT; a total of 735 participants demonstrated implicit

bias approximately 1SD above the mean ($0.9<Z_{IAT}<1.1$) on the sexuality attitudes IAT; and

1,227 participants demonstrated implicit bias approximately 1SD above the mean ($0.8<Z_{IAT}<1.2$)

on the gender stereotypes IAT.

**Table 4.**

*Model fit / selection indices and parameter estimates for +1SD data.*

|  | $\chi^2$ | $w$ | MDL | AC1 | AC2 | D | G | OB |
|---|---|---|---|---|---|---|---|---|
| Black/ |  |  |  |  |  |  |  |  |
| White |  |  |  |  |  |  |  |  |
| 1 | 38.84 | 0.03 | 15,939.90 | 0.12 (.00) | 0.08 (.00) | 0.90 (.00) | 0.56 (.06) | 0.77 (.03) |
| 2 | 425.76 | 0.08 | 16,132.80 | 0.33 (e) | 0.13 (e) | 0.84 (.00) | 0.51 (.01) | 1.00 (e) |
| 2* | 460.09 | 0.09 | 16,147.51 | 0.51 (.02) | 0.30 (.02) | 0.84 (.00) | 0.51 (.01) | Fixed |
| 2† | 425.76 | 0.08 | 16,130.35 | 0.52 (.02) | 0.31 (.02) | 0.84 (.00) | 0.51 (.06) | N/A |
| 3 | 406.66 | 0.08 | 16,122.94 | 0.27 (.08) | 0.15 (.04) | 0.88 (.01) | 0.53 (.01) | 0.76 (.14) |
| Gay/ |  |  |  |  |  |  |  |  |
| Straight |  |  |  |  |  |  |  |  |
| 1 | 100.62 | 0.04 | 20,532.94 | 0.08 (.00) | 0.06 (.00) | 0.84 (.00) | 0.51 (.01) | 0.03 (.07) |
| 2 | 115.29 | 0.04 | 20,539.71 | 0.36 (.02) | 0.26 (.02) | 0.78 (.00) | 0.49 (.01) | 0.99 (.02) |
| 2* | 115.85 | 0.04 | 20,537.54 | 0.36 (.02) | 0.26 (.02) | 0.78 (.00) | 0.49 (.01) | Fixed |
| 2† | 115.88 | 0.04 | 20,537.56 | 0.36 (.02) | 0.27 (.02) | 0.78 (.00) | 0.49 (.01) | N/A |
| 3 | 101.22 | 0.04 | 20,532.37 | 0.08 (.07) | 0.06 (.05) | 0.84 (.01) | 0.51 (.01) | 0.00 (1.03) |
| Gender/ |  |  |  |  |  |  |  |  |
| Career |  |  |  |  |  |  |  |  |
| 1 | 124.51 | 0.05 | 16,427.97 | 0.06 (.00) | 0.11 (.00) | 0.88 (.00) | 0.49 (.01) | 0.92 (.03) |
| 2 | 559.77 | 0.1 | 16,645.04 | 0.23 (e) | 0.38 (e) | 0.84 (.00) | 0.54 (.01) | 1.00 (e) |
| 2* | 586.44 | 0.1 | 16,655.92 | 0.22 (.02) | 0.37 (.02) | 0.84 (.00) | 0.54 (.01) | Fixed |
| 2† | 559.77 | 0.1 | 16,642.59 | 0.23 (.02) | 0.38 (.02) | 0.84 (.00) | 0.54 (.01) | N/A |
| 3 | 517.34 | 0.09 | 16,623.52 | 0.04 (e) | 0.07 (e) | 0.88 (.00) | 0.52 (.01) | 0.00 (e) |
| 3‡ | 517.37 | 0.09 | 16,621.60 | 0.04 (.00) | 0.07 (.00) | 0.88 (.00) | 0.52 (.01) | Fixed |

*Note*: AC1 refers to associations between the (normatively) higher status group and the
(normatively) more positive attribute, and AC2 refers to associations between the (normatively)
lower status group and the (normatively) more negative attribute.  For the Black/White IAT,
AC1 represents White-pleasant associations and AC2 represents Black-unpleasant associations.

For the gay/straight IAT, AC1 represents straight-pleasant associations and AC2 represents gay-unpleasant associations. For the gender/career IAT, AC1 represents male-career associations and AC2 represents female-family associations. *OB fixed @ 0.99 , †OB deleted, ‡OB fixed @ 0.01. (Standard Errors). (e)=unable to estimate standard error.

For the racial attitudes IAT, Model 1 ($\chi^2$=38.84, MDL=15,939.90) provided better fit by wide margins across model selection indices than Model 2 ($\chi^2$=425.76, MDL=16,132.80) and Model 3($\chi^2$=406.66, MDL=16,122.94). Several issues arose when Model 2 was fit to these data: the model was unable to estimate standard errors for AC1 (i.e., White-pleasant associations), AC2 (i.e., Black-unpleasant associations), and OB, and estimated the OB parameter at ceiling. Fixing the OB parameter at 0.99 ($\chi^2$=460.09, MDL=16,147.51), as well as deleting it from the model entirely ($\chi^2$= 425.76, MDL=16,130.35), allowed standard errors to be estimated, but did not improve model fit.

For the sexuality attitudes IAT, Model 1 ($\chi^2$=100.62, MDL=20,532.94) approximately equivalent fit to Model 3 ($\chi^2$=101.22, MDL=20,532.37), and both provided better fit than Model 2 ($\chi^2$=115.29, MDL=20,539.71). Because Model 3 estimated no contribution of OB (0.00), we fixed the OB parameter in Model 3 to 0.01 ($\chi^2$=101.23, MDL=20,530.45). This constraint did not unambiguously improve fit relative to Model 1. Because Model 2 estimated OB near ceiling (0.99), we fixed OB to 0.99 ($\chi^2$=115.85, MDL=20,537.54) and also deleted it from the model entirely ($\chi^2$=115.88, MDL=20,537.56), but neither improved model fit. Like Model 3, Model 1 estimated very little contribution of OB (0.03), so the two models collapse to essentially the same model, as they did with the full sample sexual orientation data. Thus, these data are not diagnostic for discriminating between Models 1 and 3, but can be taken as evidence in favor of these models relative to Model 2.

For the gender stereotypes IAT, Model 1 ($\chi^2$=124.51, MDL=16,427.97) provided better fit by wide margins across model selection indices than Model 2 ($\chi^2$=559.77, MDL=16,645.04) and Model 3 ($\chi^2$=517.34, MDL=16,623.52). Several issues arose when Model 2 was fit to these data: the model was unable to estimate standard errors for AC1 (i.e., male-career associations), AC2 (i.e., female-family associations), and OB, and estimated the OB parameter at ceiling. Fixing the OB parameter at 0.99 ($\chi^2$=586.44, MDL=16,655.92), as well as deleting it from the model entirely ($\chi^2$=559.77, MDL=16,642.59), allowed standard errors to be estimated, but did not improve model fit. Similarly, Model 3 was unable to estimate standard errors for AC1, AC2, and OB. Fixing the OB parameter to 0.01 ($\chi^2$=517.37, MDL=16,621.60) allowed standard errors to be estimated, but did not improve model fit.

-**1SD data.** Averaged model fit indices, model selection indices, and parameter estimates for the ten 500-person samples of data randomly selected from approximately 1SD below the mean on the IAT $d$ score for each test are presented in Table 5. A total of 19,438 participants demonstrated implicit bias approximately 1SD below the mean (-1.1<$Z_{IAT}$<-0.9) on the racial attitudes IAT; a total of 951 participants demonstrated implicit bias approximately 1SD below the mean (-1.1<$Z_{IAT}$<-0.9) on the sexuality attitudes IAT; and 912 participants demonstrated implicit bias approximately 1SD below the mean (-1.2<$Z_{IAT}$<-0.8) on the gender stereotypes IAT. In contrast to the full sample data and data selected from approximately 1SD above the mean on the IAT $d$ score, participants who scored 1SD below the mean on the IAT $d$ score demonstrated bias that was either near zero, suggesting no strong preference for either group, or below zero, suggesting a preference for the normatively lower status group. Consequently, we also fit these data to models in which the target-attribute pairings were reversed to reflect these participants' apparent response tendencies. In these reversed versions of

Models 1, 2, and 3, (listed in Table 5 as Models 4, 5, and 6, respectively), two AC parameters

were estimated: one representing the extent to which associations were activated between the

(normatively) lower status group and the (normatively) more positive attribute, and another

representing the extent to which associations were activated between the (normatively) higher

status group and the (normatively) more negative attribute. For example, based on responses to

the racial attitudes IAT, the standard models estimated White-pleasant and Black-unpleasant

associations, whereas the reversed models estimated Black-pleasant and White-unpleasant

associations.

**Table 5.**

*Model fit / selection indices and parameter estimates for -1SD data.*

| | $\chi^2$ | $w$ | MDL | AC1 | AC2 | D | G | OB |
|---|---|---|---|---|---|---|---|---|
| **Black/ White** | | | | | | | | |
| 1 | 98.4 | 0.04 | 15,953.30 | 0.04 (.00) | 0.01 (.00) | 0.86 (.00) | 0.56 (.01) | 1.00 (.13) |
| 2 | 200.06 | 0.06 | 16,003.57 | 0.12 (.01) | 0.00 (.01) | 0.85 (.00) | 0.53 (.01) | 1.00 (e) |
| 2* | 204.01 | 0.06 | 16,003.09 | 0.11 (.02) | 0.00 (.02) | 0.85 (.00) | 0.53 (.01) | Fixed |
| 2† | 200.06 | 0.06 | 16,001.12 | 0.12 (.02) | 0.00 (.02) | 0.85 (.00) | 0.53 (.01) | N/A |
| 3 | 196.69 | 0.06 | 16,001.58 | 0.07 (e) | 0.00 (e) | 0.85 (.00) | 0.53 (.01) | 0.57 (e) |
| 4 | 231.07 | 0.06 | 16,019.64 | 0.01 (.00) | 0.00 (e) | 0.85 (.01) | 0.52 (.01) | 1.00 (e) |
| 4* | 231.43 | 0.06 | 16,017.22 | 0.01 (.00) | 0.00 (.00) | 0.85 (.01) | 0.52 (.01) | Fixed |
| 5 | 245.37 | 0.06 | 16,026.22 | 0.00 (.01) | 0.00 (e) | 0.85 (.01) | 0.53 (.01) | 1.00 (e) |
| 5* | 245.39 | 0.06 | 16,023.79 | 0.00 (.00) | 0.00 (.00) | 0.85 (.01) | 0.53 (.01) | Fixed |
| 6 | 245.32 | 0.06 | 16,025.90 | 0.00 (.00) | 0.00 (e) | 0.85 (.01) | 0.52 (.01) | 0.80 (.00) |
| **Gay/ Straight** | | | | | | | | |
| 1 | 107.45 | 0.04 | 21,306.71 | 0.00 (.00) | 0.02 (.00) | 0.78 (.00) | 0.53 (.01) | 0.00 (.26) |
| 2 | 56.63 | 0.03 | 21,280.74 | 0.00 (.00) | 0.05 (.01) | 0.78 (.00) | 0.51 (.01) | 0.00 (.29) |
| 2* | 56.77 | 0.03 | 21,318.91 | 0.00 (.00) | 0.05 (.01) | 0.78 (.00) | 0.51 (.01) | Fixed |
| 2† | 140.16 | 0.05 | 21,320.05 | 0.00 (.02) | 0.07 (.02) | 0.78 (.00) | 0.53 (.01) | N/A |
| 3 | 107.45 | 0.04 | 21,305.84 | 0.00 (e) | 0.02 (e) | 0.78 (.00) | 0.53 (.01) | 0.00 (e) |
| 3‡ | 107.51 | 0.04 | 21,303.94 | 0.00 (.00) | 0.02 (.00) | 0.78 (.00) | 0.53 (.01) | Fixed |
| 4 | 156.74 | 0.05 | 21,331.35 | 0.00 (e) | 0.00 (.00) | 0.77 (.00) | 0.54 (.00) | 0.10 (.31) |
| 5 | 148.07 | 0.04 | 21,326.46 | 0.01 (.01) | 0.00 (.00) | 0.77 (.01) | 0.55 (.01) | 0.20 (.42) |
| 6 | 156.71 | 0.05 | 21,330.47 | 0.00 (e) | 0.00 (.00) | 0.77 (.01) | 0.54 (.01) | 0.70 (.48) |

Gender/
Career

| Model | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| 1 | 140.67 | 0.05 | 16,705.51 | 0.01 (.00) | 0.05 (.00) | 0.85 (.00) | 0.48 (.01) | 1.00 (.08) |
| 2 | 262.68 | 0.07 | 16,765.96 | 0.00 (.01) | 0.15 (.01) | 0.84 (.00) | 0.52 (.01) | 1.00 (e) |
| 2* | 265.93 | 0.07 | 16,765.13 | 0.00 (.02) | 0.15 (.02) | 0.84 (.00) | 0.52 (.01) | Fixed |
| 2† | 262.68 | 0.07 | 16,763.51 | 0.00 (.02) | 0.15 (.02) | 0.84 (.00) | 0.52 (.01) | N/A |
| 3 | 240.01 | 0.06 | 16,754.32 | 0.00 (.00) | 0.03 (e) | 0.85 (.00) | 0.51 (.01) | 0.00 (e) |
| 3‡ | 240.02 | 0.06 | 16,752.39 | 0.00 (.00) | 0.03 (.00) | 0.85 (.00) | 0.51 (.01) | Fixed |
| 4 | 323.35 | 0.07 | 16,796.86 | 0.01 (.00) | 0.00 (.00) | 0.84 (.00) | 0.52 (.01) | 0.67 (.21) |
| 5 | 328.57 | 0.07 | 16,798.90 | 0.00 (.00) | 0.00 (.00) | 0.84 (.00) | 0.52 (.01) | 0.20 (.42) |
| 6 | 324.33 | 0.07 | 16,796.48 | 0.01 (.00) | 0.00 (e) | 0.84 (.00) | 0.52 (.01) | 0.00 (e) |
| 6‡ | 324.35 | 0.07 | 16,794.56 | 0.01 (.00) | 0.00 (.00) | 0.84 (.00) | 0.52 (.01) | Fixed |

*Note*: For Models 1-3, AC1 refers to associations between the (normatively) higher status group and the (normatively) more positive attribute, and AC2 refers to associations between the (normatively) lower status group and the (normatively) more negative attribute. For the Black/White IAT, AC1 represents White-pleasant associations and AC2 represents Black-unpleasant associations. For the gay/straight IAT, AC1 represents straight-pleasant associations and AC2 represents gay-unpleasant associations. For the gender/career IAT, AC1 represents male-career associations and AC2 represents female-family associations. For Models 4-6, these pairings were reversed, and AC1 refers to associations between the (normatively) lower status group and the (normatively) more positive attribute, and AC2 refers to associations between the (normatively) higher status group and the (normatively) more negative attribute. For the Black/White IAT, AC1 represents Black-pleasant associations and AC2 represents White-unpleasant associations. For the gay/straight IAT, AC1 represents gay-pleasant associations and AC2 represents straight-unpleasant associations. For the gender/career IAT, AC1 represents female-career associations and AC2 represents male-family associations.

*OB fixed @ 0.99 , †OB deleted, ‡OB fixed @ 0.01. (Standard Errors). (e)=unable to estimate standard error.

For the racial attitudes IAT, Model 1 ($\chi^2$=98.4, MDL=15,953.30) provided better fit by wide margins across model selection indices than Model 2 ($\chi^2$=200.06, MDL=16,003.57) and Model 3 ($\chi^2$=196.69, MDL=16,001.58), as well as the reversed Model 4 ($\chi^2$=231.07, MDL=16,019.64), Model 5 ($\chi^2$=245.37, MDL=16,026.22), and Model 6 ($\chi^2$=245.32, MDL=16,025.90). Model 2 was unable to estimate standard error for OB, and estimated the OB parameter at ceiling. Fixing the OB parameter at 0.99 ($\chi^2$=204.01, MDL=16,003.09), as well as deleting it from the model entirely ($\chi^2$=200.06, MDL=16,001.12), allowed standard errors to be estimated, but did not improve model fit. Model 3 was unable to estimate standard errors for AC1 (i.e., White-pleasant associations), AC2 (i.e., Black-unpleasant associations), and OB.

Fixing the OB parameter to 0.01 ($\chi^2$=200.87, MDL=16,001.74) resolved the standard error issue, but did not improve model fit. Model 4 was unable to estimate standard errors for AC2 (i.e., White-unpleasant associations) and OB. Fixing the OB parameter to 0.99 ($\chi^2$=231.43, MDL=16,017.22) resolved the standard error issues, but did not significantly improve model fit. Similarly, Model 5 was unable to estimate standard errors for AC2 and OB. Fixing the OB parameter to 0.99 ($\chi^2$=245.39, MDL=16,023.79) resolved the standard error issues, but only slightly improved model fit. Model 6 was unable to estimate standard errors for AC4.

For the sexuality attitudes IAT data, Model 2 ($\chi^2$=56.63, MDL=21,280.74) provided better fit by wide margins across model selection indices than Model 1 ($\chi^2$=107.45, MDL=21,306.71) and Model 3 ($\chi^2$=107.45, MDL=21,305.84), as well as the reversed Model 4 ($\chi^2$=156.74, MDL=21,331.35), Model 5 ($\chi^2$=148.07, MDL= 21,326.46), and Model 6 ($\chi^2$= 156.71, MDL= 21,330.47). Model 2 estimated no contribution of OB, and fixing OB to 0.01 ($\chi^2$=56.77, MDL=21,318.91) did not improve model fit, whereas deleting OB from the model ($\chi^2$=140.16, MDL=21,320.05) significantly decreased model fit. Model 3 was unable to estimate standard errors for AC1 (i.e., straight-pleasant associations), AC2 (i.e., gay-unpleasant associations), and OB. Fixing the OB parameter to 0.01 ($\chi^2$=107.51, MDL=21,303.94) allowed standard errors to be estimated, but did not improve model fit. Models 4 and 6 were unable to estimate standard errors for AC2 (i.e., straight-unpleasant associations).

For the gender stereotypes IAT data, Model 1 ($\chi^2$=140.67, MDL=16,705.51) provided better fit by wide margins across model selection indices than Model 2 ($\chi^2$=262.68, MDL=16,765.96) and Model 3 ($\chi^2$=240.01, MDL=16,754.32), as well as the reversed Model 4($\chi^2$=323.35, MDL=16,796.86), Model 5 ($\chi^2$=328.57, MDL=16,798.90), and Model 6 ($\chi^2$=324.33, MDL=16,796.48).  Model 2 was unable to estimate standard error for OB, and

estimated the OB parameter at ceiling. Fixing the OB parameter at 0.99 ($\chi^2$=265.93,

MDL=16,765.13), as well as deleting it from the model entirely ($\chi^2$=262.68, MDL=16,763.51),

allowed standard errors to be estimated, but did not improve model fit. Similarly, Model 3 was

unable to estimate standard errors for AC2 (i.e., female-family associations) and OB. Fixing the

OB parameter to 0.01 ($\chi^2$=240.02, MDL=16,752.39) allowed standard errors to be estimated, but

did not improve model fit. Model 6 was unable to estimate standard errors for AC2 (i.e., male-

family associations) and OB. Fixing the OB parameter to 0.01 ($\chi^2$= 324.35, MDL= 16,794.56)

resolved the standard error issues, and only marginally improved model fit.

### General Discussion

The Quadruple process (Quad) model [7] specifies the influence of four qualitatively

distinct processes to responses on implicit measures. The model has been fit to a wide variety of

empirical data (Table 1), and its stochastic and construct validity have been established [7, 8].

The purpose of the present analyses was to compare three versions of the Quad model, each

based on different theoretical assumptions, in order to determine which provides best fit to

empirical data. Model 1 (i.e., the standard version of the model) assumes that the inhibitory

process (OB) only influences responses to target trials (e.g., pictures of Black and White people)

in the incompatible block of the IAT (e.g., when White and unpleasant share a response key),

and that the accuracy-oriented detection (D) process can only drives responses on incompatible

attribute trials when biased associations (AC) are not activated.  Model 2 retains Model 1's

assumptions about when OB has influence, but assumes that biased associations (AC) can only

drive responses on incompatible attribute trials when detection (D) fails.  Model 3 assumes that

OB drives responses on both incompatible target and attribute trials when AC and D would

produce conflicting responses and, thus, obviates the issue of AC versus D dominance on

attribute trials. Models 4-6 retain the assumptions of Models 1-3, respectively, but reverse the target-attribute pairings and were only applied to data from participants who demonstrated non-normative response tendencies.

Each of these versions of the Quad model was fit to randomly-selected subsets of real participants' data from very large sets of IAT data spanning three content domains. Additionally, in order to test the robustness of the assumptions of each model, they were fit to data drawn from one standard deviation above and below the mean on an index of bias for each domain.

Model l, the standard version of the Quad model, fit data from the racial attitudes IAT better than did the other specifications of the model. This was true for randomly-selected data, as well as for data drawn from participants who scored either 1 SD above or below the mean on the IAT $d$ score. In other words, Model 1 appears to fit data well from participants with a wide range of racial biases. Regardless of participants' level of bias, Model 1 estimated strong influence of accuracy-oriented detection (D) and inhibitory (OB) processes, as well as a general positivity bias (G). Additionally, Model 1 estimated stronger AC1 (i.e., White-pleasant) than AC2 (i.e., Black-unpleasant) associations, suggesting that racial bias as measured by this IAT is driven more strongly by pro-White than anti-Black associations, which is congruent with previous research in the intergroup bias literature (e.g., [27]). Interestingly, this was even true for participants with bias 1SD below the mean. These participants' average IAT $d$ score ($M$=-0.09, $SD$=0.02) indicates a slight evaluative preference for Blacks over Whites, seemingly in contrast to their positive White-pleasant and Black-unpleasant associations. Given that the Quad model estimates parameters based on response accuracy whereas the IAT $d$ score is based on response latency, this may simply reflect a gap between accuracy- and latency-based measures of bias. Alternately, this pattern of results may reflect one limitation of the Quad model: it cannot

estimate associations between two different groups and the same evaluative dimension (e.g., White-pleasant, Black-pleasant). Thus, though the standard sets of associations estimated by the Quad model (e.g., White-pleasant, Black-unpleasant) are congruent with most participants' response tendencies [19], this overlooks associations (e.g., White-unpleasant, Black-pleasant) that may influence participants with the opposite response tendencies. However, all of the reversed models estimated only weak AC1 (i.e., Black-pleasant associations) and no influence of AC2 (i.e., White-unpleasant associations). Moreover, none of the reversed models fit data from participants with racial bias 1SD below the mean better than the standard models. Thus, the standard version of the Quad model provided best fit for these participants, though more research on participants with unusually low levels of bias is warranted. Taken together, these data suggest that, in the domain of racial attitudes, the assumptions that the inhibitory process (OB) only influences responses to target trials (i.e., pictures of Black and White people) in the incompatible block of the IAT (i.e., when White and unpleasant share a response key) and that the accuracy-oriented detection (D) process can only drive responses on incompatible attribute trials when biased associations (AC) are not activated provide the best fit to data.

Model l also fit data from the gender stereotypes IAT better than did the other specifications of the model. This was true for randomly-selected data, as well as data drawn from participants who scored either 1 SD above or below the mean on the IAT *d* score. As it did with the racial attitudes IAT, Model 1 appears to fit data well from participants with a wide range of gender stereotypes. Regardless of participants' level of bias, Model 1 estimated strong influence of accuracy-oriented detection (D) and inhibitory (OB) processes, as well as a tendency to make family-related responses in the absence of other guides (G). Additionally, Model 1 estimated stronger AC2 (i.e., female-family) than AC1 (i.e., male-career) associations. This pattern of

results adds nuance to the interpretation of the average IAT *d* scores on the gender stereotypes IAT, which were greater than zero for all three subsets of data (random, +1SD, -1SD). Positive *d* scores in this case indicate stronger associations between males and careers and between females with families than between males and families and between females with careers. Given that it is a relative measure, the *d* score in and of itself does not indicate whether the association between males and careers is stronger, weaker, or equal to the association between females and families; instead, a positive *d* score only indicates that males are associated with careers more strongly than they are with families, and/or that females are associated with families more strongly than they are with careers. However, the results of these Quad model analyses indicate that female-family associations exert more influence on responses than male-career associations. Though speculative, this could be the result of more women entering the workforce, thereby reducing associations between careers and males, yet at the same time men not becoming increasingly associated with the family domain. Taken together, these data suggest that, in the domain of gender stereotypes, the assumptions that the inhibitory process (OB) only influences responses to target trials (e.g., male and female names) in the incompatible block of the IAT (i.e., when male and family share a response key) and that the accuracy-oriented detection (D) process can only drives responses on incompatible attribute trials when biased associations (AC) are not activated provide the best fit to data.

In contrast to the racial attitudes and gender stereotypes IATs, Model 1 did not unambiguously provide best fit to data from the sexuality attitudes IAT. Though there are myriad possible differences between the sexuality IAT and the other two IATs used in the present analyses, it is clear that this IAT was more difficult for participants to do than the other two IATs: participants made incorrect responses on an average of 11.74% of the sexuality IAT trials,

in contrast to 7.81% and 8.02% incorrect responses on the racial attitudes and gender stereotypes

IATs, respectively. Additionally, the low-bias (-1SD) sample of participants on the sexuality IAT

also had lowest average IAT *d* score (*M*=-0.14) compared to the low-bias participants on the

racial attitudes IAT (*M*=-0.09) *t*(998)=20.76, *p*<.0001, *d*=1.31, and gender stereotypes IAT

(M=0.03) t(998)=59.37, *p*<.001, *d*=3.76, suggesting perhaps greater conflict or variability in

evaluations of sexuality relative to racial attitudes and gender stereotypes. Model 1 and Model 3

provided approximately equivalent model fit to randomly-selected data and data from

participants with bias 1SD above the mean, estimating strong influence of accuracy-oriented

detection (D) as well as a general positivity bias (G). Both models also estimated stronger AC1

(i.e., straight-pleasant) than AC2 (i.e., gay-unpleasant) associations, conceptually replicating the

pattern of results found with the racial attitudes IAT, and suggesting that sexuality bias as

measured by this IAT is driven more strongly by pro-straight than anti-gay associations.

However, in contrast to the racial attitudes and gender stereotypes IATs, both Model 1 and

Model 3 estimated little or no influence of inhibitory (OB) processes on responses. When OB is

small or zero, Models 1 and 3 collapse to essentially the same model, which explains why they

provide equivalent fit in these cases. In contrast, Model 2 provided best fit to data from

participants with bias 1SD below the mean. Like Model 1 and Model 3, Model 2 estimated

strong influence of accuracy-oriented detection (D), no influence of inhibitory (OB) processes,

and a very slight positivity bias (G). Model 2 estimated stronger AC2 (i.e., gay-unpleasant) than

AC1 (i.e., straight-pleasant) associations for participants with bias 1SD below the mean, which is

surprising given that the average IAT *d* score for these participants (*M*=-0.14, *SD*=0.05) indicates

an evaluative preference for gay over straight people. Once again, this may reflect a gap between

accuracy- and latency-based measures of bias. Alternately, given that most people demonstrate

an implicit evaluative preference for straight over gay people [19], it may also reflect a limitation

of the Quad model to estimate associations for people with unusually low levels of bias.

However, all of the reversed models estimated only weak AC1 (i.e., gay-pleasant) associations

and no influence of AC2 (i.e., straight-unpleasant) associations. Moreover, none of the reversed

models fit data from participants with sexual orientation bias 1SD below the mean better than the

standard models. Taken together, these data suggest that inhibitory (OB) processes exert little to

no influence on an implicit measure of sexuality attitudes. Moreover, these data suggest that

activated associations (AC) only drive responses when accuracy-oriented detection (D) fails for

participants who are relatively low in implicit sexuality bias.

**Implications and Limitations**

The results of the present analyses indicate that the standard version of the Quad model

(Model 1) provides best fit for most of the data tested here: Model 1 outperformed Model 2 and

Model 3 on all subsets of the racial attitudes and gender stereotypes IATs, as well as reversed

Models 4-6 for participants with non-normative response tendencies. These findings are

encouraging, given that Model 1 is the version of the model that has been used in almost all

research utilizing the Quad model published to date. However, these results also indicate that that

model specification is not a one-size-fits-all endeavor. In the domain of sexuality attitudes, the

data do not allow us to discriminate between Model 1 and Model 3 in terms of best fit for

randomly-selected participants and participants with relatively strong pro-straight bias. In

contrast, Model 2 provided superior fit to data from participants with relatively pro-gay attitudes.

Because Model 2 assumes that biased associations (AC) can only drive responses on

incompatible trials when accuracy-oriented detection (D) fails, this may indicate that people with

relatively pro-gay attitudes have only weakly biased associations to begin with. Alternately, it

may reflect changing norms at the population level: both implicit and explicit attitudes towards gay people relative to straight people grew more positive between 2006-2013 [28], whereas implicit and explicit attitudes towards Black people relative to White people remained virtually unchanged during roughly the same period [29].

One implication of these findings is that different assumptions about the way in which cognitive processes interact to drive responses on the IAT may be appropriate for different populations and content domains.  Thus, researchers may need to tailor specific versions of the Quad model to specific tasks, situations, or populations of interest. For example, Gonsalkorale, von Hippel, Sherman, & Klauer [30] applied a modified version of the Quad model to data from the Go/No-Go Association Task (GNAT, [5]). In their version of the GNAT, one target group (i.e., Muslims) was always designated a "go" item. Because the target group "Muslim" had functional priority in this task, they specified a version of the Quad model to estimate one D parameter for Muslim targets and another D parameter for all other targets. This modified version of the Quad model provided superior fit compared to the standard version of the model. Thus, the present analyses may help to inform researchers in selecting or creating the appropriate version of the model to fit their specific needs.

The analyses presented here are potentially limited in some ways. For example, they are based on respondents to the Project Implicit IAT demonstration website and should not be mistaken as representing a definable population [19]. There are a variety of ways in which self-selection influences data collected through Project Implicit: for example, in learning about and choosing to visit the site and in choice and completion of measures. Nonetheless, the size and diversity of these data are much greater than is available in most laboratory studies. Moreover, the Quad model has been successfully applied to data from a wide variety of populations,

ranging from college undergraduates on four different continents, to internet samples, to clinical populations. Thus, these data are useful and appropriate for examining the psychometric properties of the Quad model.

Another potential limitation of the present analyses is that they are based solely on empirical data, whereas model validation studies often also include simulated data. However, we feel that it is not necessary to use simulated data in the present analyses because previous research has already demonstrated the construct and predictive validity of the Quad model with a wide range of data, populations, and content domains. Additionally, Burke [31] demonstrated that the Quad model provides good fit to simulated data. The intent of the present analyses is to determine which among three theoretically-defensible versions of the model provides best fit to real data. Simulated data would need to be simulated from one version of the model, which would necessarily bias it against other versions of the model. Thus, the empirical data used in the present analyses are better suited to our model comparison goals.

**Conclusion**

The present research makes several contributions. Methodologically, it empirically demonstrates that the standard specification of the Quad model provides best fit to most data from a variety of content domains and magnitudes of bias. However, it also demonstrates that alternate specifications of the model may be appropriate for some content domains or populations. Theoretically, it helps to clarify the relative contributions of in-/majority-group favoritism and out-/minority-group derogation to intergroup bias in the domains of race and sexual orientation. It also highlights the need for additional research on participants with unusually low levels of bias. More generally, this research highlights an important advantage of a

modeling approach; namely, to empirically compare the predictive ability of competing

theoretical specificiations.

**References**

1. Gawronski B. Ten frequently asked questions about implicit measures and their frequently supposed, but not entirely correct answers. Canadian Psychology/Psychologie canadienne. 2009 Aug;50(3):141.

2. Fazio RH, Jackson JR, Dunton BC, Williams CJ. Variability in automatic activation as an unobtrusive measure of racial attitudes: a bona fide pipeline?. Journal of personality and social psychology. 1995 Dec;69(6):1013.

3. Wittenbrink B, Judd CM, Park B. Evidence for racial prejudice at the implicit level and its relationship with questionnaire measures. Journal of personality and social psychology. 1997 Feb;72(2):262.

4. Greenwald AG, McGhee DE, Schwartz JL. Measuring individual differences in implicit cognition: the implicit association test. Journal of personality and social psychology. 1998 Jun;74(6):1464.

5. Nosek BA, Banaji MR. The go/no-go association task. Social cognition. 2001 Dec 1;19(6):625-66.

6. Fazio RH, Towles-Schwen T. The MODE model of attitude-behavior processes. Dual process theories in social psychology. 1999:97-116.

7. Conrey FR, Sherman JW, Gawronski B, Hugenberg K, Groom CJ. Separating multiple processes in implicit social cognition: the quad model of implicit task performance. Journal of personality and social psychology. 2005 Oct;89(4):469.

8. Sherman JW, Gawronski B, Gonsalkorale K, Hugenberg K, Allen TJ, Groom CJ. The self-regulation of automatic associations and behavioral impulses. Psychological review. 2008 Apr;115(2):314.

9. Krieglmeyer R, Sherman JW. Disentangling stereotype activation and stereotype application in the stereotype misperception task. Journal of personality and social psychology. 2012 Aug;103(2):205.

10. Meissner F, Rothermund K. Estimating the contributions of associations and recoding in the Implicit Association Test: The ReAL model for the IAT. Journal of Personality and Social Psychology. 2013 Jan;104(1):45.

11. Nadarevic L, Erdfelder E. Cognitive processes in implicit attitude tasks: An experimental validation of the Trip Model. European Journal of Social Psychology. 2011 Mar 1;41(2):254-68.

12. Payne BK. Prejudice and perception: the role of automatic and controlled processes in misperceiving a weapon. Journal of personality and social psychology. 2001 Aug;81(2):181.

13. Payne BK, Hall DL, Cameron CD, Bishara AJ. A process model of affect misattribution. Personality and Social Psychology Bulletin. 2010 Sep 13.

14. Stahl C, Degner J. Assessing automatic activation of valence: A multinomial model of EAST performance. Experimental Psychology. 2007 Jan;54(2):99-112.

15. De Houwer J. A structural analysis of indirect measures of attitudes. The psychology of evaluation: Affective processes in cognition and emotion. 2003 Jan 30:219-44.

16. Kornblum S, Hasbroucq T, Osman A. Dimensional overlap: cognitive basis for stimulus-response compatibility--a model and taxonomy. Psychological review. 1990 Apr;97(2):253.

17. Batchelder WH, Riefer DM. Theoretical and empirical review of multinomial process tree modeling. Psychonomic Bulletin & Review. 1999 Mar 1;6(1):57-86.

18. Riefer DM, Batchelder WH. Multinomial modeling and the measurement of cognitive processes. Psychological Review. 1988 Jul;95(3):318.

19. Nosek BA, Smyth FL, Hansen JJ, Devos T, Lindner NM, Ranganath KA, Smith CT, Olson KR, Chugh D, Greenwald AG, Banaji MR. Pervasiveness and correlates of implicit attitudes and stereotypes. European Review of Social Psychology. 2007 Jan 1;18(1):36-88.

20. Lindsay DS, Jacoby LL. Stroop process dissociations: The relationship between facilitation and interference. Journal of Experimental Psychology: Human Perception and Performance. 1994 Apr;20(2):219.

21. Jacoby LL. A process dissociation framework: Separating automatic from intentional uses of memory. Journal of memory and language. 1991 Oct 31;30(5):513-41.

22. Greenwald AG, Nosek BA, Banaji MR. Understanding and using the implicit association test: I. An improved scoring algorithm. Journal of personality and social psychology. 2003 Aug;85(2):197.

23. Xu K, Nosek B, Greenwald A. Psychology data from the race Implicit Association Test on the Project Implicit demo website. Journal of Open Psychology Data. 2014 Mar 18;2(1).

24. Nosek, B.A., Banaji, M. and Greenwald, A.G., 2002. Harvesting implicit group attitudes and beliefs from a demonstration web site. Group Dynamics: Theory, Research, and Practice, 6(1), p.101.

25. Cohen J. A power primer. Psychological bulletin. 1992 Jul;112(1):155.

26. Wu H, Myung JI, Batchelder WH. Minimum description length model selection of multinomial processing tree models. Psychonomic bulletin & review. 2010 Jun 1;17(3):275-86.

27. Brewer MB. In-group bias in the minimal intergroup situation: A cognitive-motivational analysis. Psychological bulletin. 1979 Mar;86(2):307.

28. Westgate E, Riskind R, Nosek B. Implicit preferences for straight people over lesbian women and gay men weakened from 2006 to 2013. Collabra. 2015 Jul 23;1(1).

29. Schmidt K, Nosek BA. Implicit (and explicit) racial attitudes barely changed during Barack Obama's presidential campaign and early presidency. Journal of Experimental Social Psychology. 2010 Mar 31;46(2):308-14.

30. Gonsalkorale K, von Hippel W, Sherman JW, Klauer KC. Bias and regulation of bias in intergroup interactions: Implicit attitudes toward Muslims and interaction quality. Journal of Experimental Social Psychology. 2009 Jan 31;45(1):161-6.

31. Burke CT. Process dissociation models in racial bias research: Updating the analytic method and integrating with signal detection approaches. Group Processes & Intergroup Relations. 2015 May 1;18(3):402-34.

32. Cohen, J. A power primer. Psychological Bulletin. 1992 Jul; 112(1), 155-159.

**Table 1 References**

1. Jin Z, Rivers AM, Sherman JW, Chen R. Measures of Implicit Gender Attitudes May Exaggerate Differences in Underlying Associations among Chinese Urban and Rural Women. Psychologica Belgica. 2016 Jan 29;56(1).

2. Scroggins WA, Mackie DM, Allen TJ, Sherman JW. Reducing Prejudice With Labels Shared Group Memberships Attenuate Implicit Bias and Expand Implicit Group Boundaries. Personality and Social Psychology Bulletin. 2015 Dec 14:0146167215621048.

3. Huntsinger JR, Sinclair S, Kenrick AC, Ray C. Affiliative social tuning reduces the activation of prejudice. Group Processes & Intergroup Relations. 2015 May 12:1368430215583518.

4. Burke CT. Process dissociation models in racial bias research: Updating the analytic method and integrating with signal detection approaches. Group Processes & Intergroup Relations. 2015 May 1;18(3):402-34.

5. Ramos MR, Barreto M, Ellemers N, Moya M, Ferreira L, Calanchini J. Exposure to sexism can decrease implicit gender stereotype bias. European Journal of Social Psychology. 2015 Jan 1.

6. Zestcott CA, Bean MG, Stone J. Evidence of negative implicit attitudes toward individuals with a tattoo near the face. Group Processes & Intergroup Relations. 2015 Sep 21:1368430215603459.

7. Ruiz JG, Andrade AD, Anam R, Taldone S, Karanam C, Hogue C, Mintzer MJ. Group-Based Differences in Anti-Aging Bias Among Medical Students. Gerontology & geriatrics education. 2015 Jan 2;36(1):58-78.

8. Calanchini J, Sherman JW, Klauer KC, Lai CK. Attitudinal and non-attitudinal components of IAT performance. Personality and Social Psychology Bulletin. 2014 Jul 1:0146167214540723.

9. Clerkin EM, Fisher CR, Sherman JW, Teachman BA. Applying the Quadruple Process model to evaluate change in implicit attitudinal responses during therapy for panic disorder. Behaviour research and therapy. 2014 Jan 31;52:17-25.

10. Gonsalkorale K, Sherman JW, Klauer KC. Measures of Implic it Attitudes May Conceal Differences in Implicit Associations The Case of Antiaging Bias. Social Psychological and Personality Science. 2013 Aug 5:1948550613499239.

11. Calanchini J, Gonsalkorale K, Sherman JW, Klauer KC. Counter-prejudicial training reduces activation of biased associations and enhances response monitoring. European Journal of Social Psychology. 2013 Aug 1;43(5):321-5.

12. Soderberg CK, Sherman JW. No face is an island: How implicit bias operates in social scenes. Journal of Experimental Social Psychology. 2013 Mar 31;49(2):307-13.

13. O'Connor RM, Lopez-Vergara HI, Colder CR. Implicit cognition and substance use: The role of controlled and automatic processes in children. Journal of studies on alcohol and drugs. 2012 Jan;73(1):134-43.

14. Allen TJ, Sherman JW. Ego Threat and Intergroup Bias A Test of Motivated-Activation Versus Self-Regulatory Accounts. Psychological science. 2011 Mar 1;22(3):331-3.

15. Gonsalkorale K, Sherman JW, Allen TJ, Klauer KC, Amodio DM. Accounting for Successful Control of Implicit Racial Bias The Roles of Association Activation, Response Monitoring, and Overcoming Bias. Personality and Social Psychology Bulletin. 2011 Nov 1;37(11):1534-45.

16. Allen TJ, Sherman JW, Klauer KC. Social context and the self-regulation of implicit bias. Group Processes & Intergroup Relations. 2010 Mar 1;13(2):137-49.

17. Gonsalkorale K, Allen TJ, Sherman JW, Klauer KC. Mechanisms of group membership and exemplar exposure effects on implicit attitudes. Social Psychology. 2010 Aug 11.

18. Gonsalkorale K, Sherman JW, Klauer KC. Aging and prejudice: Diminished regulation of automatic race bias among older adults. Journal of Experimental Social Psychology. 2009a Feb 28;45(2):410-4.

19. Gonsalkorale K, von Hippel W, Sherman JW, Klauer KC. Bias and regulation of bias in intergroup interactions: Implicit attitudes toward Muslims and interaction quality. Journal of Experimental Social Psychology. 2009 Jan 31;45(1):161-6.

20. Beer JS, Stallen M, Lombardo MV, Gonsalkorale K, Cunningham WA, Sherman JW. The Quadruple Process model approach to examining the neural underpinnings of prejudice. NeuroImage. 2008 Dec 31;43(4):775-83.

21. Bishara AJ, Payne BK. Multinomial process tree models of control and automaticity in weapon misidentification. Journal of Experimental Social Psychology. 2009 May 31;45(3):524-34.

22. Conrey FR, Sherman JW, Gawronski B, Hugenberg K, Groom CJ. Separating multiple processes in implicit social cognition: the quad model of implicit task performance. Journal of personality and social psychology. 2005 Oct;89(4):469.

**Appendix**

Model equations. Note: All equations are organized in terms of the Black/White IAT. For the gay/straight and gender/career IATs, the normative majority group (i.e., straight, male) was substituted for White, the normative minority group (i.e., gay, female) was substituted for Black, and gender/career words were substituted for pleasant/unpleasant, respectively.

**Model 1**

Compatible trials (i.e., White-pleasant / Black-unpleasant)

White:      Correct = AC1+(1-AC1) D+(1-AC1)(1-D)G

            Incorrect = (1-AC1)(1-D)(1-G)

Black:      Correct = AC2+(1-AC2)D+(1-AC2)(1-D)(1-G)

            Incorrect = (1-AC2)(1-D)G

Pleasant:   Correct = AC1+(1-AC1) D+(1-AC1)(1-D)G

            Incorrect = (1-AC1)(1-D)(1-G)

Unpleasant: Correct = AC2+(1-AC2)D+(1-AC2)(1-D)(1-G)

            Incorrect = (1-AC2)(1-D)G

Incompatible trials (i.e., White-unpleasant / Black-pleasant)

White:      Correct = AC1*D*OB+(1-AC1)D+(1-AC1)(1-D)(1-G)

            Incorrect = AC1*D*(1-OB)+AC1(1-D)+(1-AC1)(1-D)G

Black:      Correct = AC2*D*OB+(1-AC2)D+(1-AC2)(1-D)G

            Incorrect = AC2*D*(1-OB)+AC2(1-D)+(1-AC2)(1-D)(1-G)

Pleasant:   Correct = (1-AC1)D+(1-AC1)(1-D)G

            Incorrect = AC1+(1-AC1)(1-D)(1-G)

Unpleasant: Correct = (1-AC2)D+(1-AC2)(1-D)(1-G)

$$\text{Incorrect} = AC2+(1-AC2)(1-D)G$$

## Model 2

Compatible trials (i.e., White-pleasant / Black-unpleasant)

White:    Correct $= AC1+(1-AC1) D+(1-AC1)(1-D)G$

Incorrect $= (1-AC1)(1-D)(1-G)$

Black:    Correct $= AC2+(1-AC2)D+(1-AC2)(1-D)(1-G)$

Incorrect $= (1-AC2)(1-D)G$

Pleasant:   Correct $= AC1+(1-AC1) D+(1-AC1)(1-D)G$

Incorrect $= (1-AC1)(1-D)(1-G)$

Unpleasant: Correct $= AC2+(1-AC2)D+(1-AC2)(1-D)(1-G)$

Incorrect $= (1-AC2)(1-D)G$

Incompatible trials (i.e., White-unpleasant / Black-pleasant)

White:    Correct $= AC1*D*OB+(1-AC1)D+(1-AC1)(1-D)(1-G)$

Incorrect $= AC1*D*(1-OB)+AC1(1-D)+(1-AC1)(1-D)G$

Black:    Correct $= AC2*D*OB+(1-AC2)D+(1-AC2)(1-D)G$

Incorrect $= AC2*D*(1-OB)+AC2(1-D)+(1-AC2)(1-D)(1-G)$

Pleasant:   Correct $= D+(1-AC1)(1-D)G$

Incorrect $= (1-D)AC1+(1-AC1)(1-D)(1-G)$

Unpleasant: Correct $= D+(1-AC2)(1-D)(1-G)$

Incorrect $= (1-D)AC2+(1-AC2)(1-D)G$


## Model 3

Compatible trials (i.e., White-pleasant / Black-unpleasant)

White:      Correct = AC1+(1-AC1) D+(1-AC1)(1-D)G

Incorrect = (1-AC1)(1-D)(1-G)

Black:      Correct = AC2+(1-AC2)D+(1-AC2)(1-D)(1-G)

Incorrect = (1-AC2)(1-D)G

Pleasant:   Correct = AC1+(1-AC1) D+(1-AC1)(1-D)G

Incorrect = (1-AC1)(1-D)(1-G)

Unpleasant: Correct = AC2+(1-AC2)D+(1-AC2)(1-D)(1-G)

Incorrect = (1-AC2)(1-D)G

Incompatible trials (i.e., White-unpleasant / Black-pleasant)

White:      Correct = AC1*D*OB+(1-AC1)D+(1-AC1)(1-D)(1-G)

Incorrect = AC1*D*(1-OB)+AC1(1-D)+(1-AC1)(1-D)G

Black:      Correct = AC2*D*OB+(1-AC2)D+(1-AC2)(1-D)G

Incorrect = AC2*D*(1-OB)+AC2(1-D)+(1-AC2)(1-D)(1-G)

Pleasant:   Correct = AC1*D*OB+(1-AC1)D+(1-AC1)(1-D)(G)

Incorrect = AC1*D*(1-OB)+AC1(1-D)+(1-AC1)(1-D)(1-G)

Unpleasant: Correct = AC2*D*OB+(1-AC2)D+(1-AC2)(1-D)(1-G)

Incorrect = AC2*D*(1-OB)+AC2(1-D)+(1-AC2)(1-D)(G)