

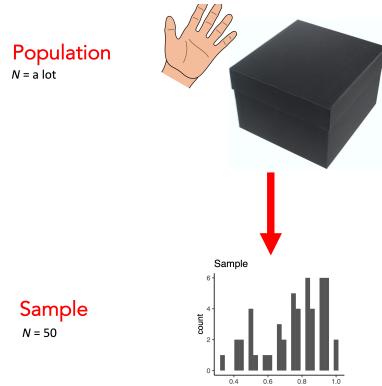
Statistical Foundations: Replication (and the failures)

2 March 2020

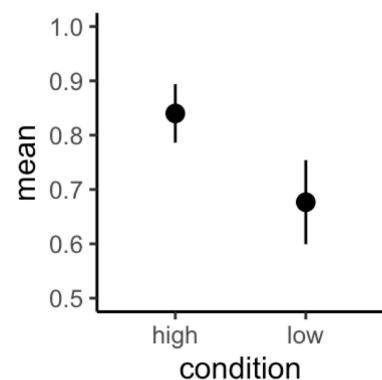
Modern Research Methods

Statistical framework for thinking about replications

Sampling



Confidence Intervals (CI)



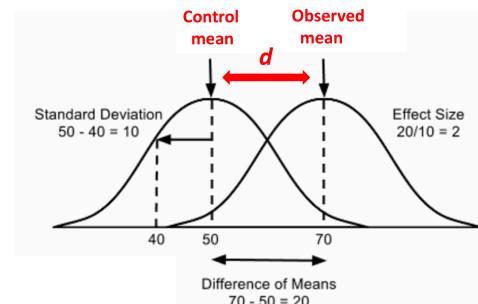
Null Hypothesis Testing

Are the means different?
(yes/no)



Effect Sizes

How different are the means?



	Original	Reproduction	Replication
Population			
Question			
Hypothesis			
Exp. Design			
Experimenter			
Data	01100 10110 11110	01100 10110 11110	01100 10110 11110
Analyst			
Code			
Estimate			
Claim			

Original



Different



REPLICATE = Get same result
with a new dataset

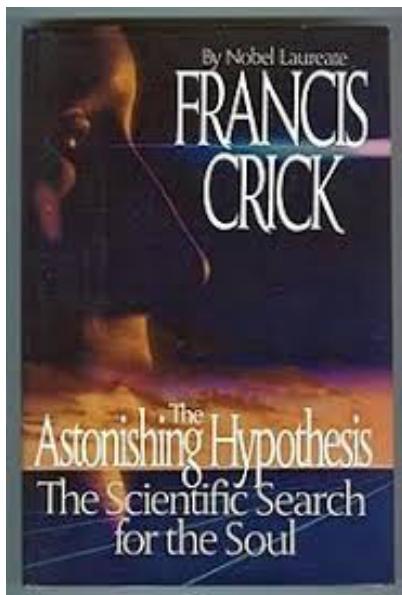
Research Article

The Value of Believing in Free Will

Encouraging a Belief in Determinism Increases Cheating

Kathleen D. Vohs¹ and Jonathan W. Schooler²

¹*Department of Marketing, Carlson School of Management, University of Minnesota, and* ²*Department of Psychology, University of British Columbia*



Read Passage

Anti-free-will essay

Consciousness essay
(control)

"glitchy" Math Test

Measure: How much did they cheat?

advocating a deterministic worldview that dismisses individual causation may similarly promote undesirable behavior. In this vein, Peale (1989) bemoaned how quickly and consistently deviant behavior is tagged a “disease,” a label that obviates personal responsibility for its occurrence. As a recent *Washington Post* article on neuroscience and moral behavior put it succinctly, “Reducing morality and immorality to brain chemistry—rather than free will—might diminish the importance of personal responsibility” (Vedantam, 2007, p. A01).

Although some people have speculated about the societal risks that might result from adopting a viewpoint that denies personal responsibility for actions, this hypothesis has not been explored empirically. In the two experiments reported here, we manipulated beliefs related to free will and measured their influence on morality as manifested in cheating behavior. We hypothesized that participants induced to believe that human behavior is under the control of predetermined forces would cheat more than would participants not led to believe that behavior is predetermined. Our experimental results supported this hypothesis.

EXPERIMENT 1

Method

Participants

Participants were 30 undergraduates (13 females, 17 males).

Procedure

Participants came to the lab individually. First, according to the condition to which they were randomly assigned, they read one of two passages from *The Astonishing Hypothesis*, a book written by Francis Crick (1994), the Nobel-prize-winning scientist. In the *anti-free-will* condition, participants read statements claiming that rational, high-minded people—including, according to Crick, most scientists—now recognize that actual free will is an illusion, and also claiming that the idea of free will is a side effect of the architecture of the mind. In the *control* condition, participants read a passage from a chapter on consciousness, which did not discuss free will. After reading their assigned material, participants completed the Free Will and Determinism scale (FWD; Paulhus & Margesson, 1994) and the Positive and Negative Affectivity Schedule (PANAS; Watson, Clark, & Tellegen, 1988), which we used to assess whether the reading manipulation affected their beliefs and mood.

Subsequently, participants were given a computer-based mental-arithmetic task (von Hippel, Lakin, & Shakarchi, 2005) in which they were asked to calculate the answers to 20 problems (e.g., $1 + 8 + 18 - 12 + 19 - 7 + 17 - 2 + 8 - 4 = ?$), presented individually. They were told that the computer had a programming glitch and the correct answer would appear on the screen while they were attempting to solve each problem, but that they could stop the answer from being displayed by pressing

the space bar after the problem appeared. Furthermore, participants were told that although the experimenter would not know whether they had pressed the space bar, they should try to solve the problems honestly, on their own. In actuality, the computer had been rigged not only to show the answers, but also to record the number of space-bar presses. The dependent measure of cheating was the number of times participants pressed the space bar to prevent the answer from appearing. Afterward, participants were debriefed and thanked for their participation.

Results

Scores on the FWD Scale

We first checked to see whether participants’ beliefs about free will were affected by the excerpts they read (anti-free-will vs. control condition). As expected, scores on the Free Will subscale of the FWD scale showed that participants in the anti-free-will condition reported weaker free-will beliefs ($M = 13.6, SD = 2.66$) than participants in the control condition ($M = 16.8, SD = 2.67$), $t(28) = 3.28, p < .01$. Scores on the other three subscales of the FWD scale (Fate, Scientific Causation, and Chance) did not differ as a function of condition, $t < 1$.

Cheating

We first recoded the dependent measure by subtracting the number of space-bar presses from 20, so that higher scores indicated more cheating. Analysis of the main dependent measure, degree of cheating, revealed that, as predicted, participants cheated more frequently after reading the anti-free-will essay ($M = 14.00, SD = 4.17$) than after reading the control essay ($M = 9.67, SD = 5.58$), $t(28) = 3.04, p < .01$.

Does Rejecting the Idea of Free Will Lead to Cheating?

To test our hypothesis that cheating would increase after participants were persuaded that free will does not exist, we first calculated the correlation between scores on the Free Will subscale and cheating behavior. As expected, we found a strong negative relationship, $r(30) = -.53$, such that weaker endorsement of the notion that personal behavior is determined by one’s own will was associated with more instances of cheating.

We next performed a mediation analysis to test our prediction that degree of belief in free will would determine degree of cheating. Using analysis of covariance (ANCOVA), we found support for this hypothesis: When Free Will subscale scores were entered as a predictor of cheating alongside experimental condition, the effect of condition failed to predict cheating, $F < 1$, whereas the effect of free-will beliefs remained significant, $F(1, 27) = 7.81, p < .01$.

Ancillary Measure: Mood

To ensure that the essays did not inadvertently alter participants’ moods, we assessed positive and negative emotions using the PANAS. Mood did not differ between conditions, $t < 1.35, ps > .19$.

Cheating

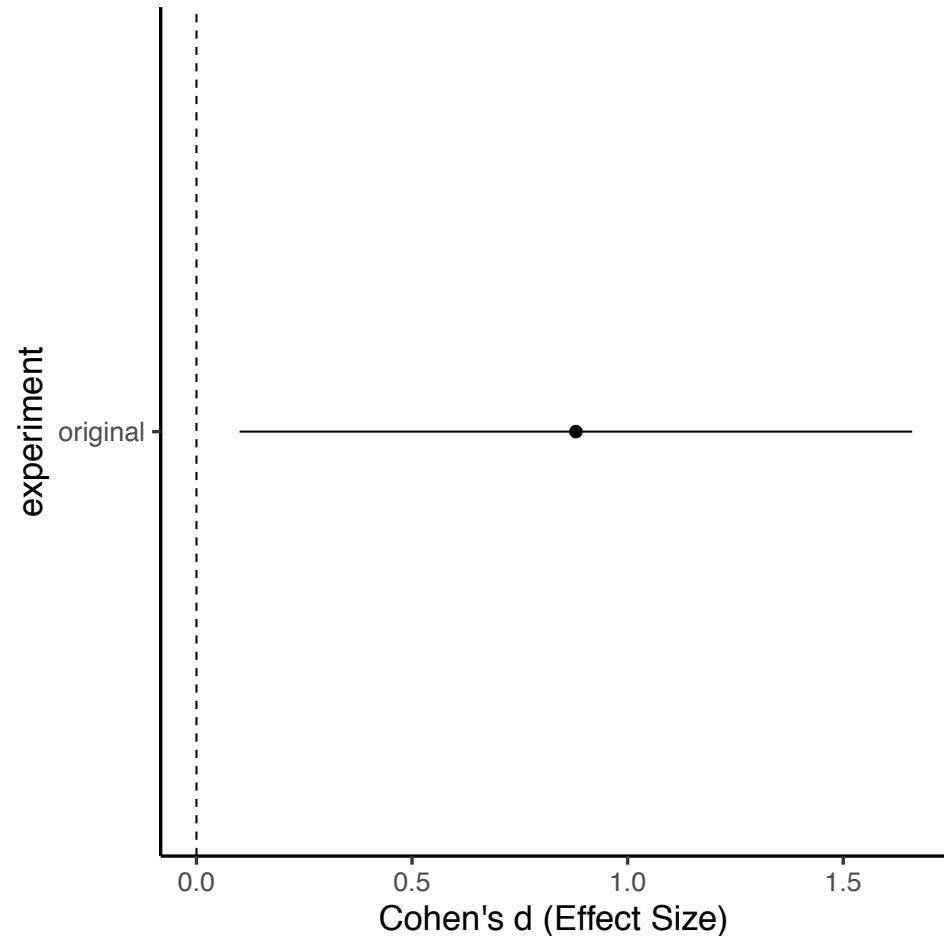
We first recoded the dependent measure by subtracting the number of space-bar presses from 20, so that higher scores indicated more cheating. Analysis of the main dependent measure, degree of cheating, revealed that, as predicted, participants cheated more frequently after reading the anti-free-will essay ($M = 14.00, SD = 4.17$) than after reading the control essay ($M = 9.67, SD = 5.58$), $t(28) = 3.04, p < .01$.

Calculating an effect size

$$\text{Effect Size} = \frac{\text{diff. between means}}{\text{standard dev.}}$$

Cheating

We first recoded the dependent measure by subtracting the number of space-bar presses from 20, so that higher scores indicated more cheating. Analysis of the main dependent measure, degree of cheating, revealed that, as predicted, participants cheated more frequently after reading the anti-free-will essay ($M = 14.00$, $SD = 4.17$) than after reading the control essay ($M = 9.67$, $SD = 5.58$), $t(28) = 3.04$, $p < .01$.



Replication of Vohs and Schooler (2017)

$N = 58$

```
> t.test(determinism, freewill)
```

Welch Two Sample t-test

```
data: determinism and freewill  
t = 0.77137, df = 50.951, p-value = 0.4441  
alternative hypothesis: true difference in means is not equal to 0
```

95 percent confidence interval:

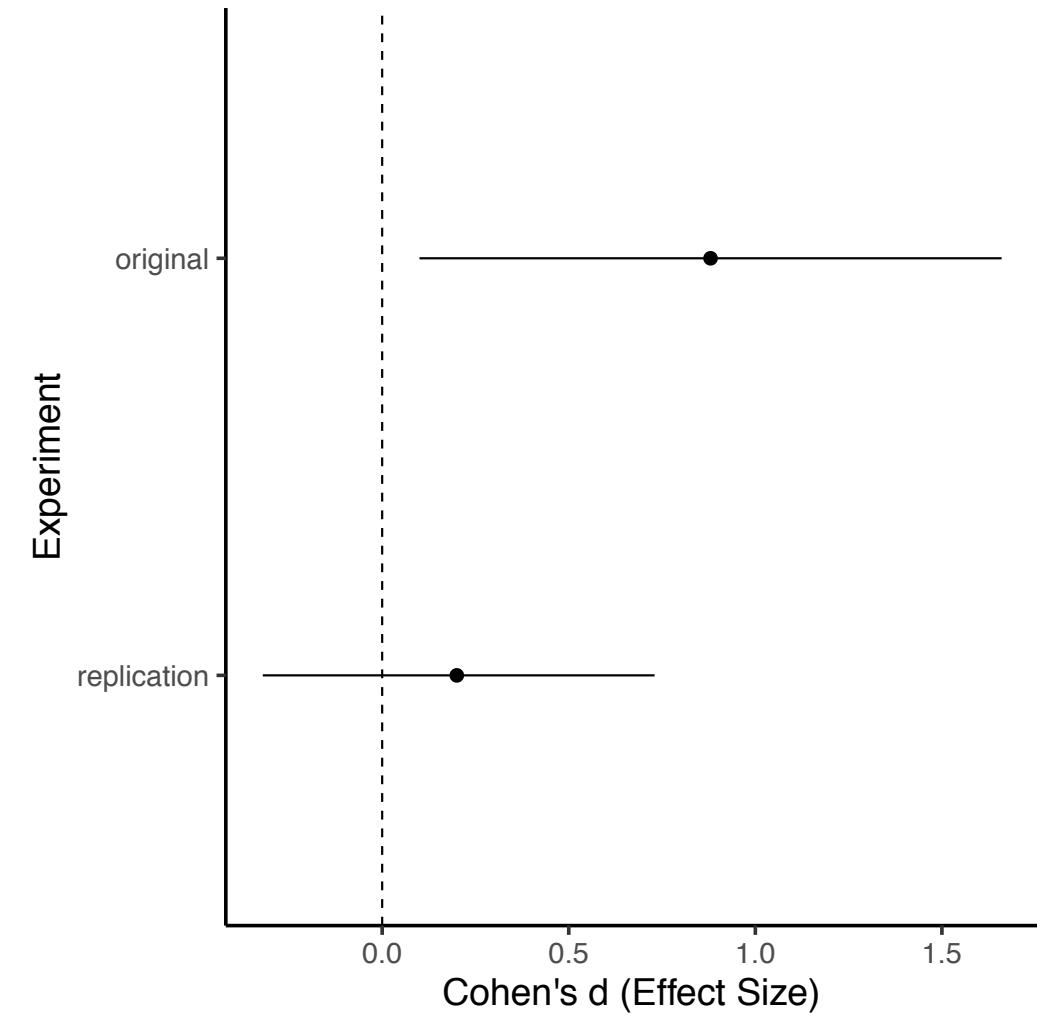
-1.934273 4.348066

sample estimates:

mean of x mean of y

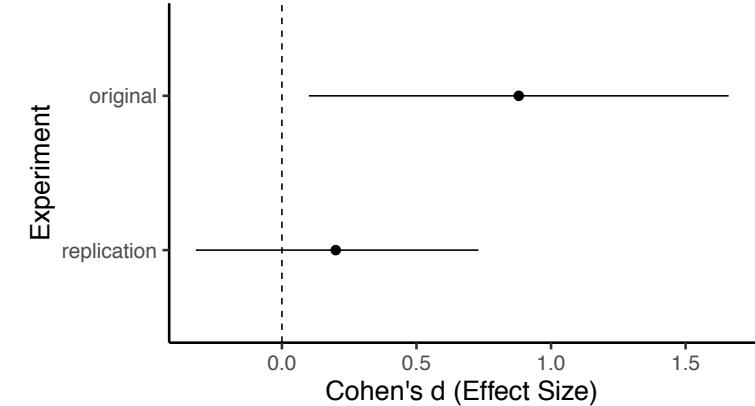
6.793103 5.586207

```
> mes(6.793103, 5.586207, 6.831541, 4.931801, 29, 29,  
+     verbose = F) %>%  
+     select(d, l.d, u.d) %>%  
+     rename(ci_lower = l.d,  
+           ci_upper = u.d)  
d ci_lower ci_upper  
1 0.2      -0.32      0.73
```



Interpreting the replication

Statistics	Original	Replication	Interpretation
p -value (t -test)	<.01	.44	Did not replicate
Effect size	.84 [.06, 1.62]	.2 [-.32, .72]	Effect size much smaller (1/4), and the confidence interval for the replicates includes 0.



Talk to the person(s) next to you, and generate a list of reasons why the Vohs & Schooler effect might not have replicated.

Why might an effect not replicate?

1. Effect isn't real, got unlucky in original.
2. Effect is real, but got unlucky in replication.

	Retain H_0	Reject H_0
H_0 is true	correct	Error (Type I)
H_0 is false	Error (Type II)	correct

How replicable are psychological studies?

Rumblings that something isn't right...

Essay

Why Most Published Research Findings Are False

John P. A. Ioannidis

Puzzlingly High Correlations in fMRI Studies of Emotion, Personality, and Social Cognition¹

Edward Vul,¹ Christine Harris,² Piotr Winkielman,² & Harold Pashler²

¹Massachusetts Institute of Technology and ²University of California, San Diego

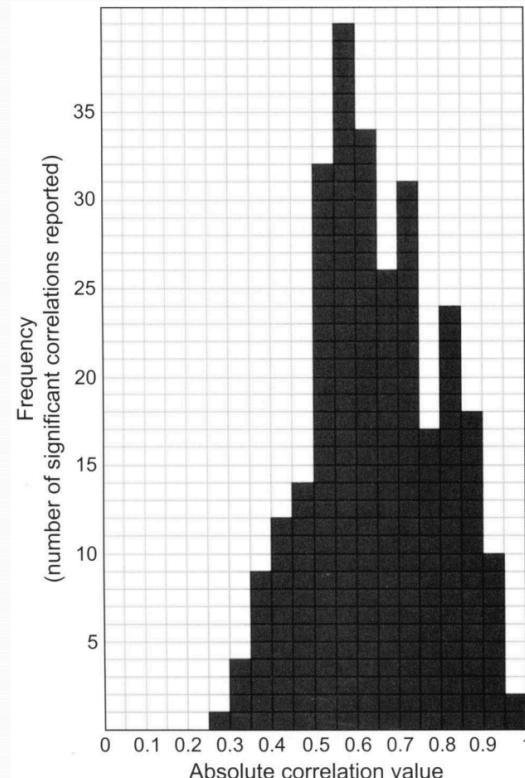


Fig. 1. A histogram of the correlations between evoked blood oxygenation level dependent response and behavioral measures of individual differences seen in the studies identified for analysis in the current article.

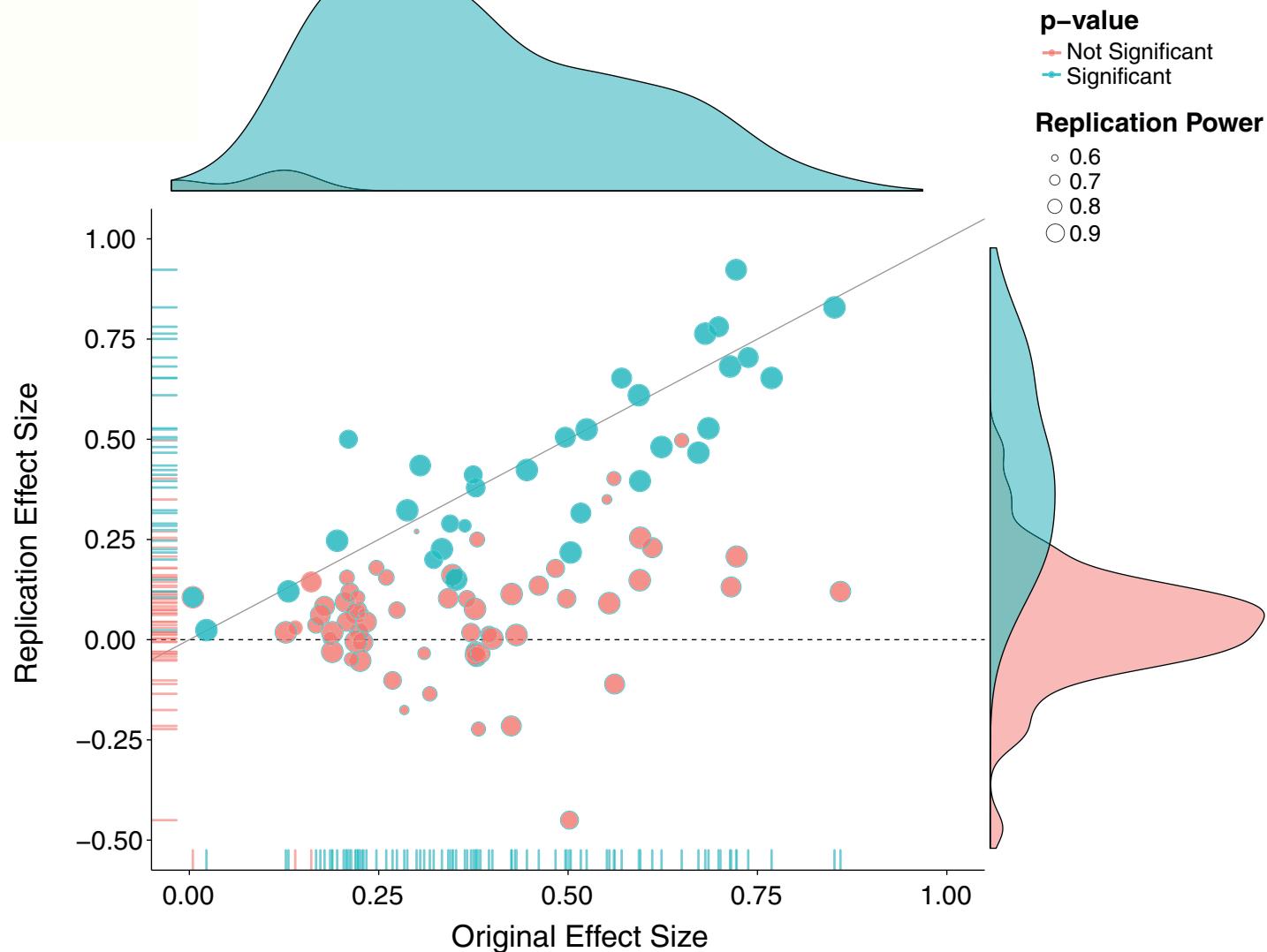
Estimating the reproducibility of psychological science

Open Science Collaboration* (2015)

Conducted replications of 100 psychology effects

Replication effect size half the size of the original, on average.

97 of original studies had $p < .05$;
only 37 of replications (47 in CI
of original)



Why might an effect not replicate?

1. Effect isn't real, got unlucky in original.
2. Effect is real, but got unlucky in replication.

	Retain H_0	Reject H_0
H_0 is true	correct	Error (Type I)
H_0 is false	Error (Type II)	correct

- More than 5% of studies are failing to replicate – why?
- There must be other reasons that researchers are getting a positive effect, when there is no effect.

Reason # 1: Fraud (i.e, researchers are just making up their data)

Harvard Dean Confirms Misconduct in Hauser Investigation

by [Greg Miller](#) on 20 August 2010, 3:11 PM | [2 Comments](#)

[Email](#) [Print](#) | [f](#) [t](#) [+1](#) [0](#)

[Reddit](#) [SU](#) [More](#)

[PREVIOUS ARTICLE](#)

[NEXT ARTICLE](#)

In an e-mail sent earlier today to Harvard Law and Sciences (FAS), confirms that cognitive scientists have been asked to undergo a thorough investigation by a faculty member "to determine whether the researcher's work under FAS standards."

November 13, 2011

Fraud Scandal Fuels Debate Over Practices of Social Psychology

Even legitimate researchers cut corners, some admit



By Christopher Shea

The discovery that the Dutch researcher Diederik A. Stapel made up the data for dozens of research papers has shaken up the field of social psychology, fueling a discussion not just about outright fraud, but also about subtler ways of misusing research data. Such misuse

Reason #2: Error in reporting/analysis

Current Biology
Retraction



Retraction Notice to: The Negative Association between Religiousness and Children's Altruism across the World

Jean Decety,* Jason M. Cowell, Kang Lee, Randa Mahasneh, Susan Malcolm-Smith, Bilge Selcuk, and Xinyue Zhou

*Correspondence: decety@uchicago.edu

<https://doi.org/10.1016/j.cub.2019.07.030>

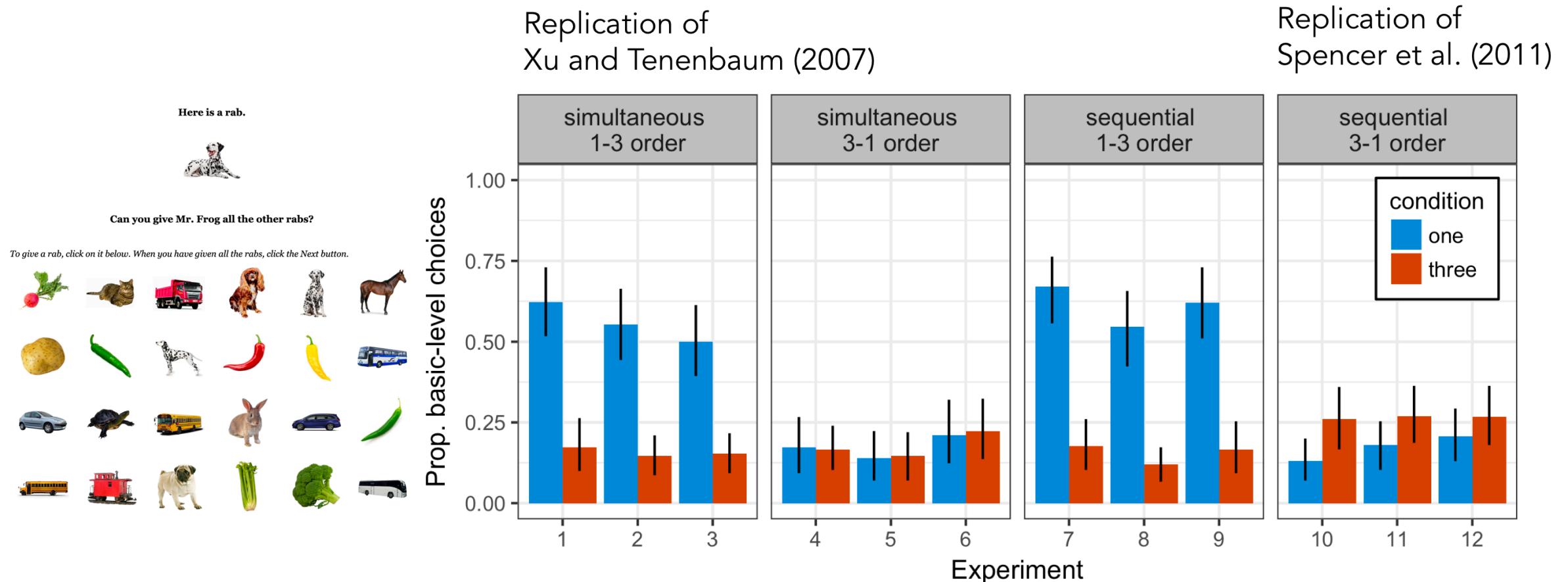
(Current Biology 25, 2951–2955; November 16, 2015)

In our paper, we reported cross-cultural differences in how the religious environment of a child negatively impacted their sharing, their judgments of the actions of others, and how their parents evaluated them. An error in this article, our incorrect inclusion of country of origin as a covariate in many analyses, was pointed out in a correspondence from Shariff, Willard, Muthukrishna, Kramer, and Henrich (<https://doi.org/10.1016/j.cub.2016.06.031>). When we reanalyzed these data to correct this error, we found that country of origin, rather than religious affiliation, is the primary predictor of several of the outcomes. While our title finding that increased household religiousness predicts less sharing in children remains significant, we feel it necessary to explicitly correct the scientific record, and we are therefore retracting the article. We apologize to the scientific community for any inconvenience caused.

Coded country as continuous measure, rather than categorical.

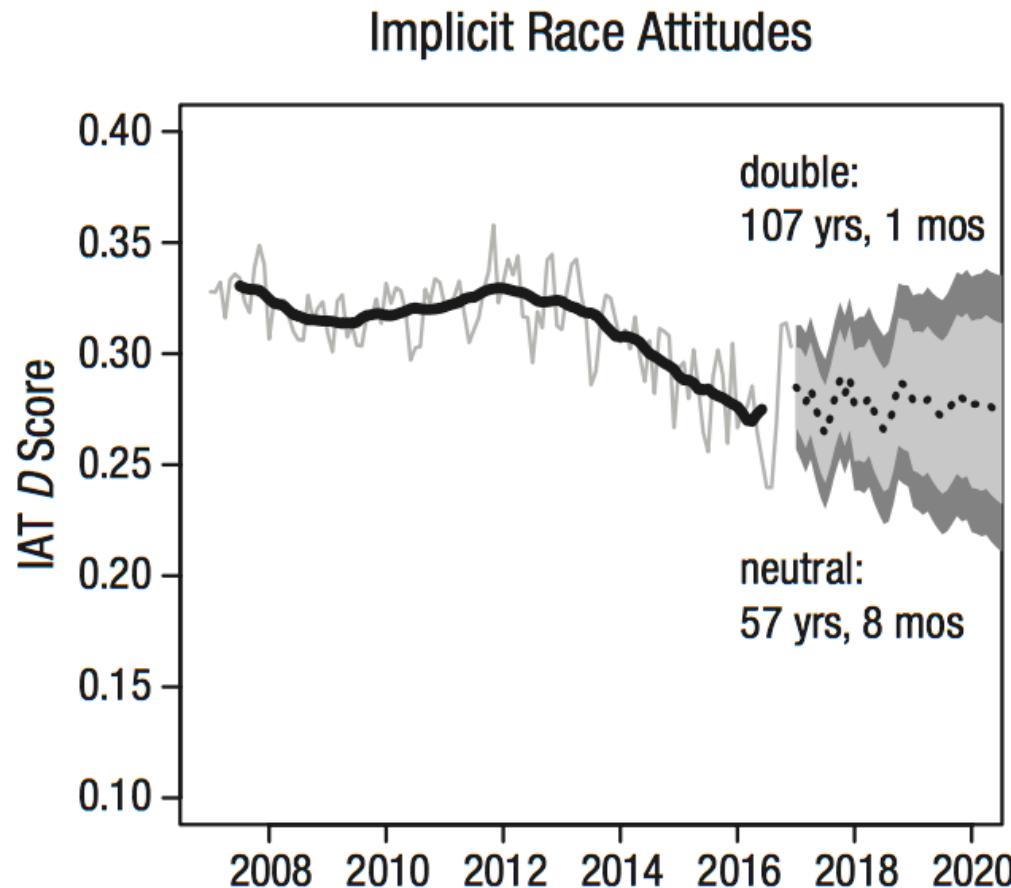
- 1 – US
- 2 – Canada
- 3 – Japan

Reason #3: Hidden moderator



(Lewis and Frank, 2016)

Reason #4: Actual change in population



(Charlesworth & Banaji, 2019)

Reason #5: File Drawer Problem (or, "publication bias")

Psychological Bulletin
1979, Vol. 86, No. 3, 638-641

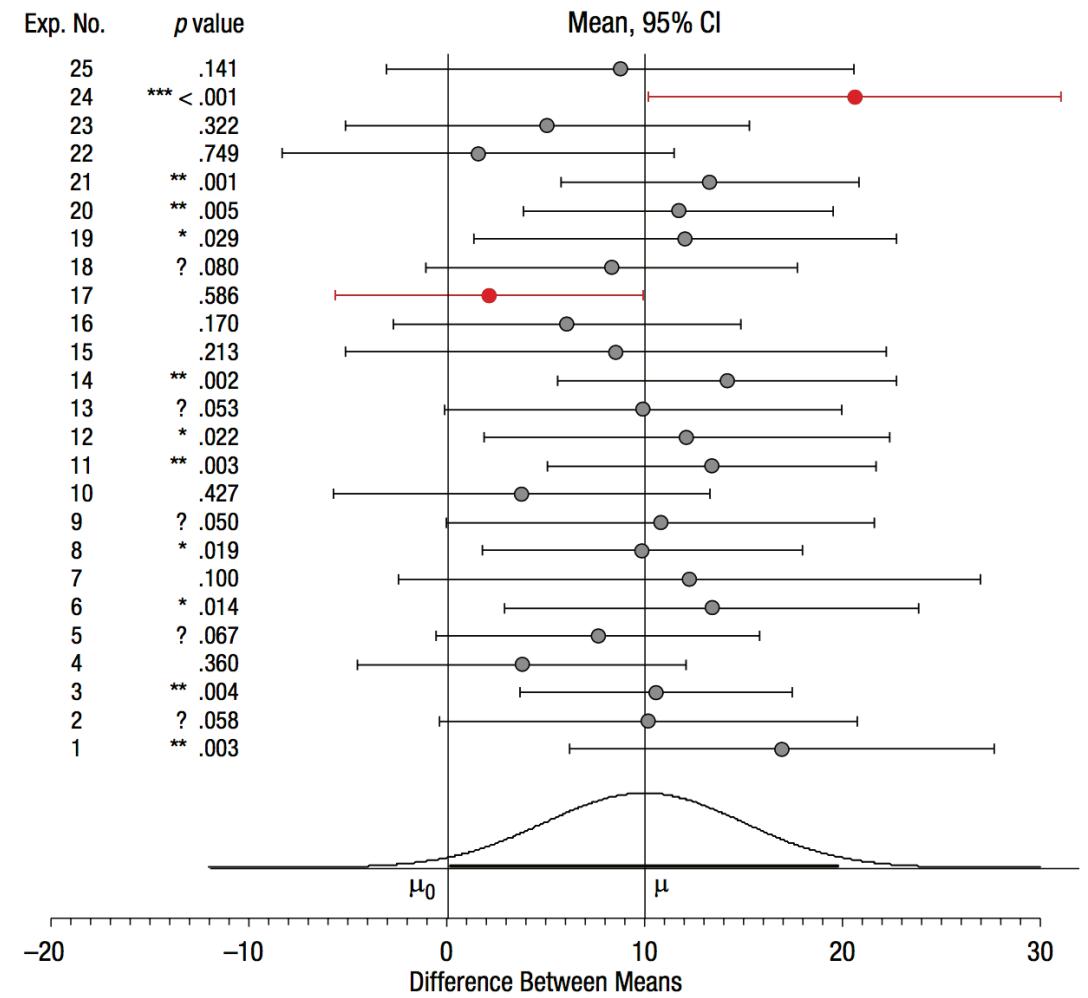
The "File Drawer Problem" and Tolerance for Null Results

Robert Rosenthal
Harvard University



Historically, researchers have only been able to publish results if $p < .05$.

But, they run many studies, and put the studies that "fail" in their file drawer.



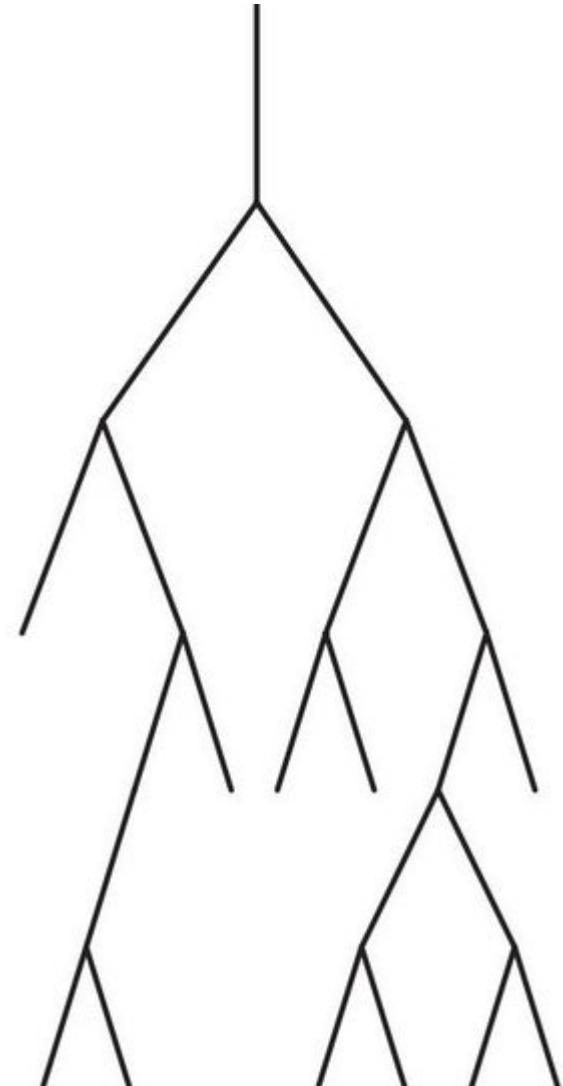
Reason #6: Data-dependent analysis

The Statistical Crisis in Science

Data-dependent analysis—a “garden of forking paths”—explains why many statistically significant comparisons don’t hold up.

Choosing your analysis based on seeing your data/the outcome of a test (“analytic flexibility”)

“p-hacking”



"...it is unacceptably easy to publish *statistically significant* evidence consistent with *any hypothesis*" -- Simmons, Nelson, & Simonsohn, 2011

"Researcher degrees of freedom" -- flexibility in data collection, analysis, and reporting

Collect more data?

Should some observations be excluded? Which ones?

Which conditions should be combined with which ones?

Which measures should we analyze? Should we transform the measure?

Which control variables should we consider?

Next Time: Solutions to the replicability crisis

Reading:

Education ► Schools Teachers Universities Students

Peer review and scientific publishing

Chris Chambers, Marcus Munafo and more than 80 signatories

Wed 5 Jun 2013
07.45 EDT

[f](#) [t](#) [m](#)

Trust in science would be improved by study pre-registration

Open letter: We must encourage scientific journals to accept studies before the results are in

