# Correcting for bias in the literature

## A comprehensive comparison of meta-analytic methods for bias-correction

Felix Schönbrodt, Evan Carter, Will Gervais, Joe Hilgard

Felix Schönbrodt
Ludwig-Maximilians-Universität
München

RESEARCH TRANSPARENCY

OSC
LMU Open Science Center

www.nicebread.de
www.researchtransparency.org
@nicebread303

LUDWIG-MAXIMILIANS-UNIVERSITÄT MÜNCHEN

Meta-analysis is at the top of the evidence-based medicine pyramid - the pinnacle of evidence-based medicine.

Cochrane Collaboration

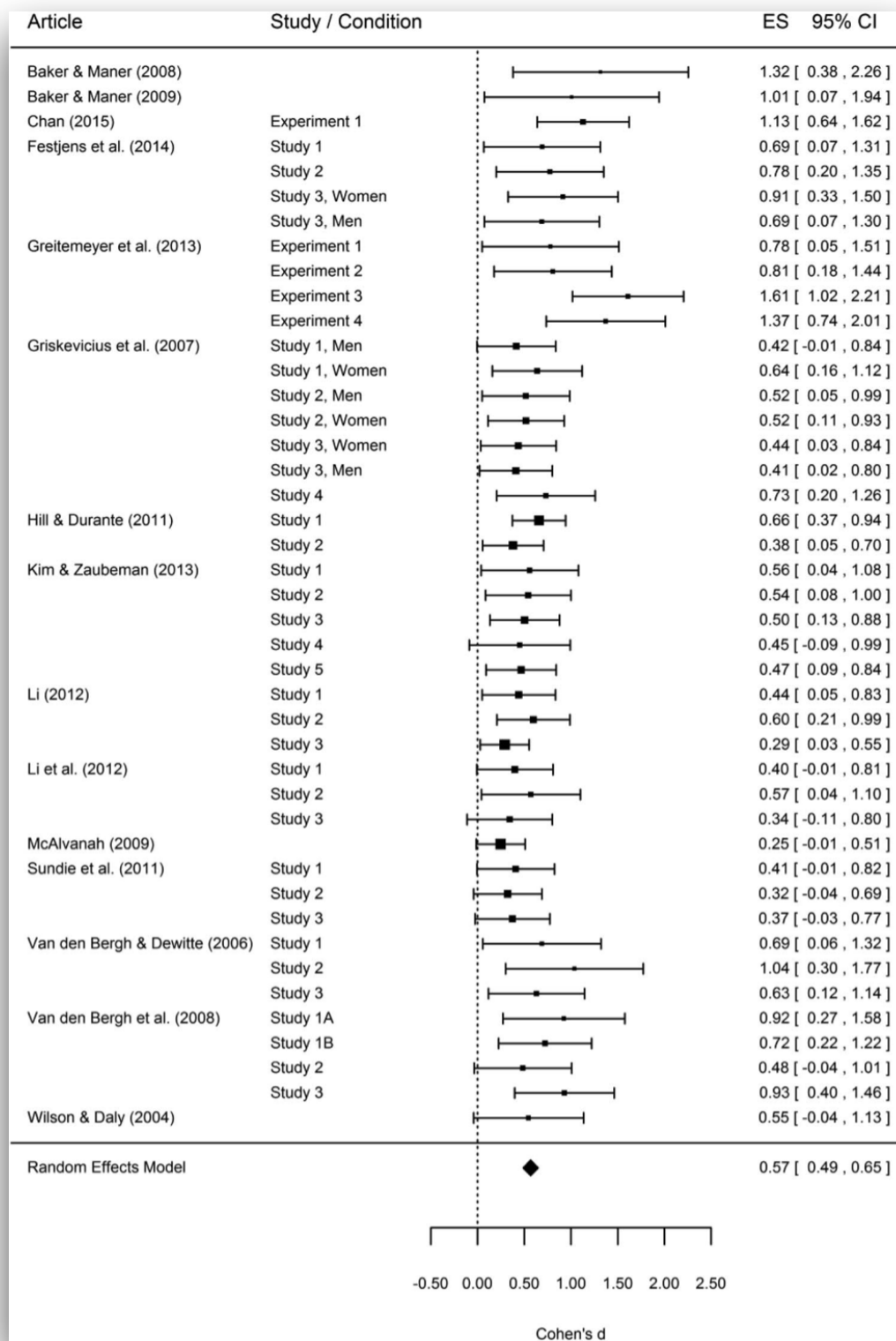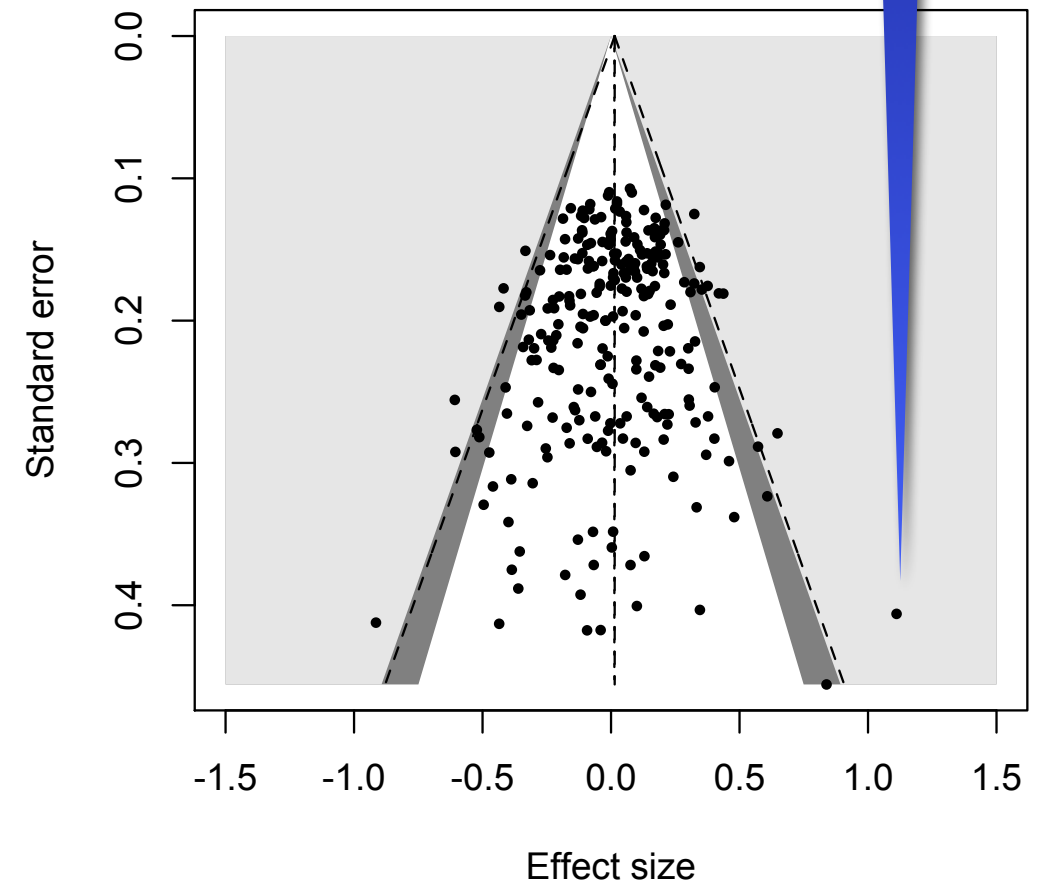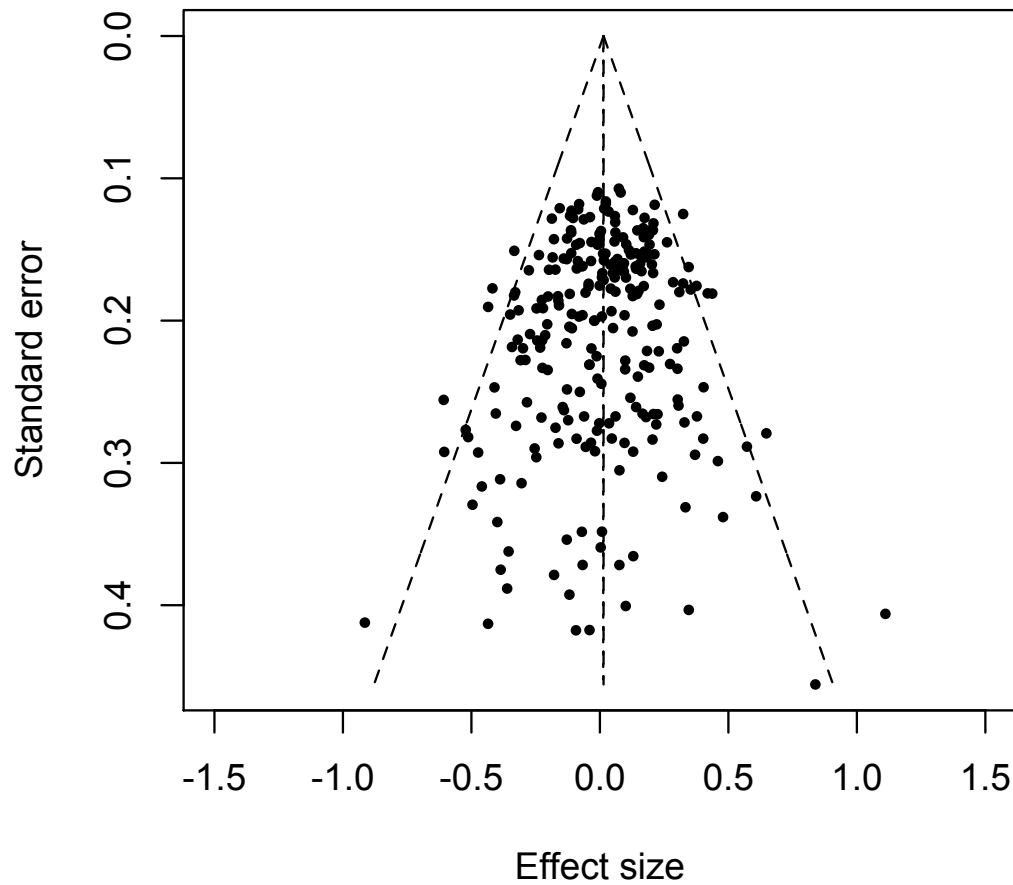Meta-analyses are fucked.

Mickey Inzlicht

| Article | Study / Condition | | ES 95% CI |
|---|---|---|---|
| Baker & Maner (2008) | | | 1.32 [ 0.38 , 2.26 ] |
| Baker & Maner (2009) | | | 1.01 [ 0.07 , 1.94 ] |
| Chan (2015) | Experiment 1 | | 1.13 [ 0.64 , 1.62 ] |
| Festjens et al. (2014) | Study 1 | | 0.69 [ 0.07 , 1.31 ] |
| | Study 2 | | 0.78 [ 0.20 , 1.35 ] |
| | Study 3, Women | | 0.91 [ 0.33 , 1.50 ] |
| | Study 3, Men | | 0.69 [ 0.07 , 1.30 ] |
| Greitemeyer et al. (2013) | Experiment 1 | | 0.78 [ 0.05 , 1.51 ] |
| | Experiment 2 | | 0.81 [ 0.18 , 1.44 ] |
| | Experiment 3 | | 1.61 [ 1.02 , 2.21 ] |
| | Experiment 4 | | 1.37 [ 0.74 , 2.01 ] |
| Griskevicius et al. (2007) | Study 1, Men | | 0.42 [ -0.01 , 0.84 ] |
| | Study 1, Women | | 0.64 [ 0.16 , 1.12 ] |
| | Study 2, Men | | 0.52 [ 0.05 , 0.99 ] |
| | Study 2, Women | | 0.52 [ 0.11 , 0.93 ] |
| | Study 3, Women | | 0.44 [ 0.03 , 0.84 ] |
| | Study 3, Men | | 0.41 [ 0.02 , 0.80 ] |
| | Study 4 | | 0.73 [ 0.20 , 1.26 ] |
| Hill & Durante (2011) | Study 1 | | 0.66 [ 0.37 , 0.94 ] |
| | Study 2 | | 0.38 [ 0.05 , 0.70 ] |
| Kim & Zaubeman (2013) | Study 1 | | 0.56 [ 0.04 , 1.08 ] |
| | Study 2 | | 0.54 [ 0.08 , 1.00 ] |
| | Study 3 | | 0.50 [ 0.13 , 0.88 ] |
| | Study 4 | | 0.45 [ -0.09 , 0.99 ] |
| | Study 5 | | 0.47 [ 0.09 , 0.84 ] |
| Li (2012) | Study 1 | | 0.44 [ 0.05 , 0.83 ] |
| | Study 2 | | 0.60 [ 0.21 , 0.99 ] |
| | Study 3 | | 0.29 [ 0.03 , 0.55 ] |
| Li et al. (2012) | Study 1 | | 0.40 [ -0.01 , 0.81 ] |
| | Study 2 | | 0.57 [ 0.04 , 1.10 ] |
| | Study 3 | | 0.34 [ -0.11 , 0.80 ] |
| McAlvanah (2009) | | | 0.25 [ -0.01 , 0.51 ] |
| Sundie et al. (2011) | Study 1 | | 0.41 [ -0.01 , 0.82 ] |
| | Study 2 | | 0.32 [ -0.04 , 0.69 ] |
| | Study 3 | | 0.37 [ -0.03 , 0.77 ] |
| Van den Bergh & Dewitte (2006) | Study 1 | | 0.69 [ 0.06 , 1.32 ] |
| | Study 2 | | 1.04 [ 0.30 , 1.77 ] |
| | Study 3 | | 0.63 [ 0.12 , 1.14 ] |
| Van den Bergh et al. (2008) | Study 1A | | 0.92 [ 0.27 , 1.58 ] |
| | Study 1B | | 0.72 [ 0.22 , 1.22 ] |
| | Study 2 | | 0.48 [ -0.04 , 1.01 ] |
| | Study 3 | | 0.93 [ 0.40 , 1.46 ] |
| Wilson & Daly (2004) | | | 0.55 [ -0.04 , 1.13 ] |
| Random Effects Model | | | 0.57 [ 0.49 , 0.65 ] |

-0.50  0.00  0.50  1.00  1.50  2.00  2.50

Cohen's d

Random effects meta-analytic estimate:
$d$ = 0.57 [0.49; 0.65]

42/43 studies are significant (98% success rate)
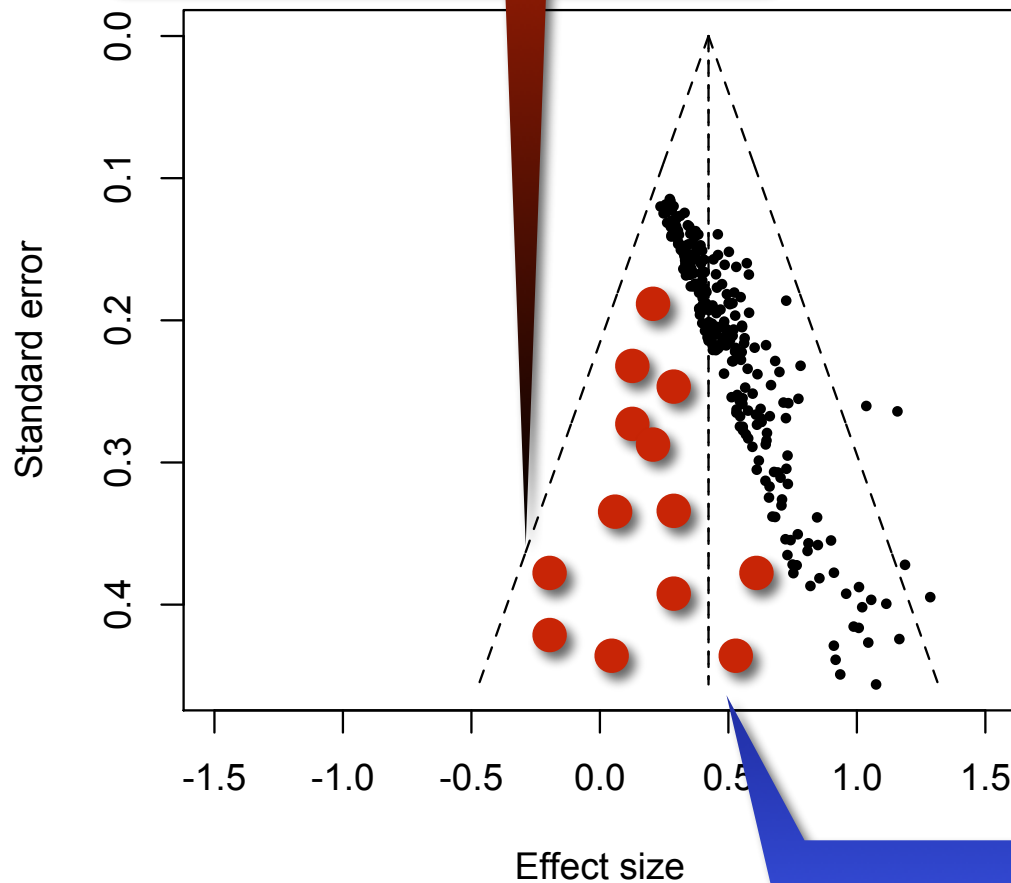
3

# True H$_0$ samples*

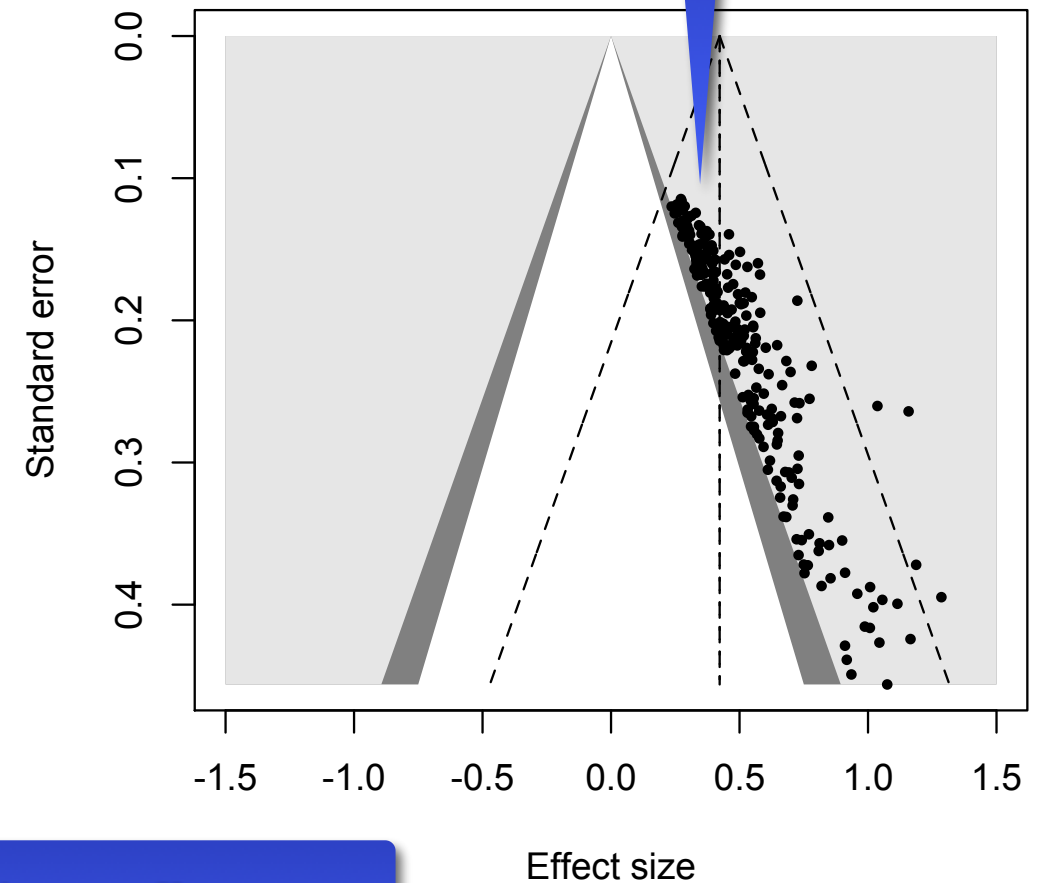5% false positive ("significant") studies

* simulated data

# True H₀ + directional publication bias



There seem to be some studies missing!

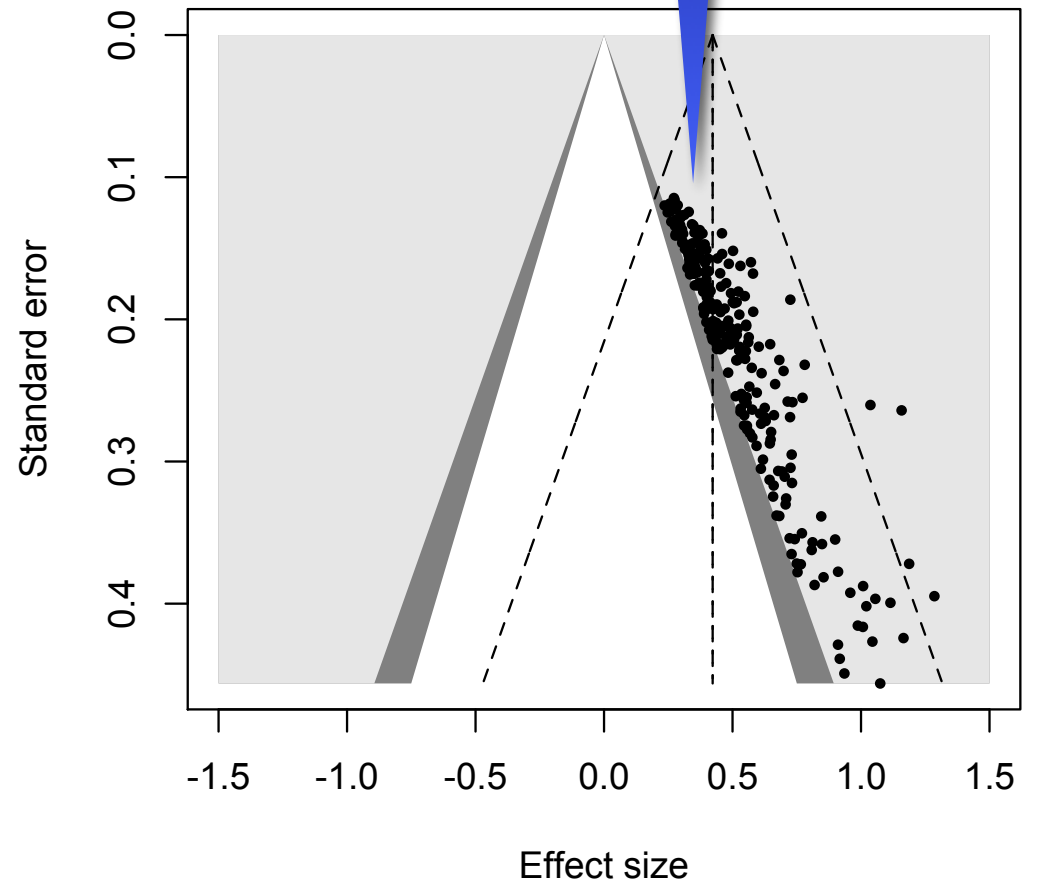Studies "huddle" against the significance threshold

Meta-analytic effect size estimate: d = 0.42

Standard error

Effect size

* simulated data

# True H₀ + publication bias

\* simulated data

# Romance, Risk, and Replication: Can Consumer Choices and Risk-Taking Be Primed by Mating Motives?

David R. Shanks
University College London

Miguel A. Vadillo
King's College London

Benjamin Riedel, Ashley Clymo, Sinita Govind, Nisha Hickin, Amanda J. F. Tamman, and Lara M. C. Puhlmann
University College London

# Romance, Risk, and Replication: Can Consumer Choices and Risk-Taking Be Primed by Mating Motives?

David R. Shanks
University College London

Miguel A. Vadillo
King's College London

Benjamin Riedel, Ashley Clymo, Sinita Govind, Nisha Hickin, Amanda J. F. Tamman, and Lara M. C. Puhlmann
University College London

14 replication studies, all *n.s.*

8

# Correcting for publication bias (PB)

*or*

# Can we clean up the mess, if we only had the right tool?

# Trim & Fill

- Originally designed as a *test* for PB, but also used to *correct* for PB

- Algorithmically fill in missing studies to achieve a symmetric funnel plot

- Compute meta-analysis on the data set including imputed studies

There seem to be some studies missing!
➔ trim-and-fill

Duval, S. & Tweedie, R. (2000). Trim and fill: a simple funnel-plot–based method of testing and adjusting for publication bias in meta-analysis. *Biometrics, 56 (2)*, 455–463.

# PET / PEESE

- Extrapolates the „small study effect" to samples with ∞ sample size

- What would be the effect size if we had an infinitely large sample?

- PET: linear regression

- PEESE: squared slope



PET (linear)    PEESE (squared)

Stanley, T. D., & Doucouliagos, H. (2013). Meta-regression approximations to reduce publication selection bias. Research Synthesis Methods, 5(1), 60–78. http://doi.org/10.1002/jrsm.1095

| |

# Selection models

- Explicitly model the functional form of publication bias

- Provide estimates for, e.g., *Prob*(published | n.s.)

- Three-parameter SM: $\mu$, $\tau$, and *Prob*(published | n.s.)

- Four-parameter SM: $\mu$, $\tau$, and *Prob*(pub | n.s. & correct direction) and *Prob*(pub | wrong direction)



McShane, B. B., Böckenholt, U., & Hansen, K. T. (2016).
Iyengar, S. & Greenhouse, J. B. (1988)
Hedges, L. V. (1984)

from Guan & Vandekerckhove, 2015)        12

# Performance of bias correcting methods

# Simulation study



Primary sample size derived from psychological literature

Table 1
*Simulation parameters*

| Experimental factors | Levels |
| --- | --- |
| True underlying effect ($\delta$) | 0, 0.2, 0.5, 0.8 |
| Between-study heterogeneity ($\tau$) | 0, 0.2, 0.4 |
| Number of studies in the meta-analytic sample ($k$) | 10, 30, 60, 100 |
| Publication bias ($PB$) | None, medium, strong |
| QRP environment ($QRP$) | None, medium, high |

fully crossed:
432 conditions

**Estimators:**
(naive) Random effects meta-analysis, Trim&Fill, PET, PEESE, PET-PEESE, three-parameter selection model (3PSM), four-parameter selection model (4PSM), $p$-curve, $p$-uniform, WAAP-WLS

14

# Results (a selection)

http://shinyapps.org/apps/metaExplorer/

# Hypothesis test

**How many % of original studies are submitted to publication bias?:**

● 0%   ○ 60%   ○ 90%

**Heterogeneity (tau):**

● 0   ○ 0.2   ○ 0.4

**Number of studies in meta-analysis:**

○ 10   ● 30   ○ 60   ○ 100

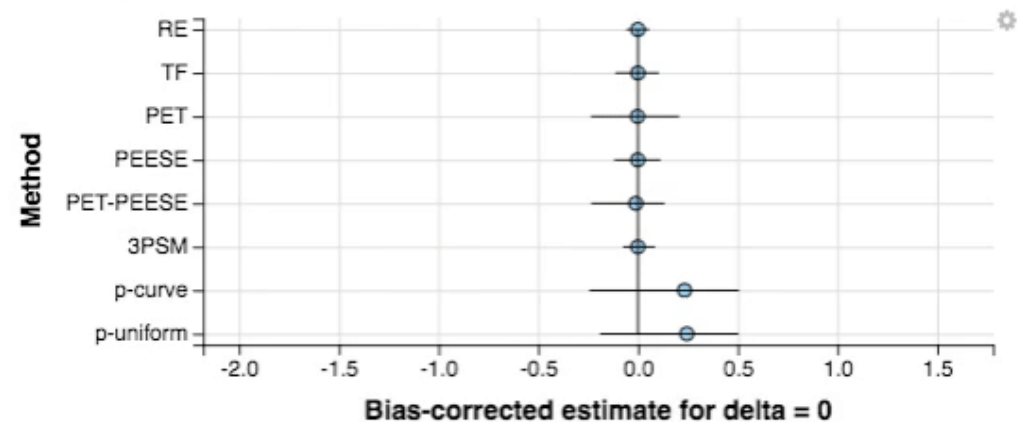**True effect size under H1 (for power computation)**

● 0.2   ○ 0.5   ○ 0.8

**QRP environment:**

● none   ○ med   ○ high

## Under H0

If in reality there is no effect: What is the probability that a method falsely concludes 'There is an effect'?



# Effect size estimation

## Basic settings

**How many % of original studies are submitted to publication bias?:**

● 0%   ○ 60%   ○ 90%

**Heterogeneity (tau):**

○ 0   ● 0.2   ○ 0.4

**Number of studies in meta-analysis:**

○ 10   ○ 30   ○ 60   ● 100

**True effect size under H1 (for power computation)**

● 0.2   ○ 0.5   ○ 0.8

**QRP environment:**

● none   ○ med   ○ high

## Bias-corrected estimates of the true effect

Under H0



17

# Method performance check

- Hope that all bias-correcting methods will converge on the same value? Usually that does not happen

- ➡No vote counting - no triangulation:

  - Even if three out of four methods converge on a value this is irrelevant, when those three are known to perform badly in plausible conditions.

- Use the app to see which bias-correcting methods perform well in plausible conditions for the meta-analysis at hand

- Do a sensitivity analysis - but only including methods that passed the performance check!

## Meta-analysis - the pinnacle of evidence-based research?

## Meta-analyses are fucked?

- Publication bias and $p$-hacking massively distorts the evidence:
  **Garbage in - garbage out.**
- Even meta-analyses of many dozen significant primary studies can come from a null effect.
- Each type of bias-correction works in some conditions, but fails in other conditions.
  *Problem*: We do not know which condition we are in.
- Doing biased research and hoping to correct it afterward *does not work*.
- Better put efforts into improving primary studies themselves (e.g., by using registered reports which combat both $p$-hacking and publication bias)

19

**Correcting for bias in psychology: A comparison of meta-analytic methods**

Evan C. Carter* — U.S. Army Research Laboratory, Aberdeen, MD, USA
Felix D. Schönbrodt* — Ludwig-Maximilians-Universität, Munich, Germany
Will M. Gervais — University of Kentucky, Lexington, KY, USA
Joseph Hilgard — University of Pennsylvania, Philadelphia, PA, USA

https://psyarxiv.com/9h3nu/

- „Researchers should **not expect** to produce a conclusive, **debate-ending result** by conducting a meta-analysis on an existing literature"

- „Instead, we imagine meta-analyses may serve best to draw attention to the existing strengths and/or weaknesses in a literature and these results can then inspire a careful re-examination of methodology and theory followed by, if necessary, **large-scale, preregistered replication efforts.**"