

# Learn to $p$ -hack like the pros!



J.J. at the English language [Wikipedia](#)

# SCIENCE

The **smartest heads** in the world immerse themselves into a research topic for years.

In that process, they become the experts – nobody knows more about that topic. The boundaries of knowledge have been **pushed forward**.

When the scientist are confident in their findings, they publish them in the best scientific journals, with the highest standards of **quality, rigor, and integrity**.

**HOW MUCH OF THAT LITERATURE  
DO YOU THINK IS TRUE?**

Researchers are not rewarded for being right,  
but rather for publishing a lot.

Nelson, Simmons, & Simonsohn (2012); Nosek, Spies, Motyl (2012); Munafò (2016)

**Shit Academics Say**  
@AcademicsSay

A. I get paid to think.  
B. About what.  
C. Tenure mostly.

RETWEETS 186 FAVORITEN 306

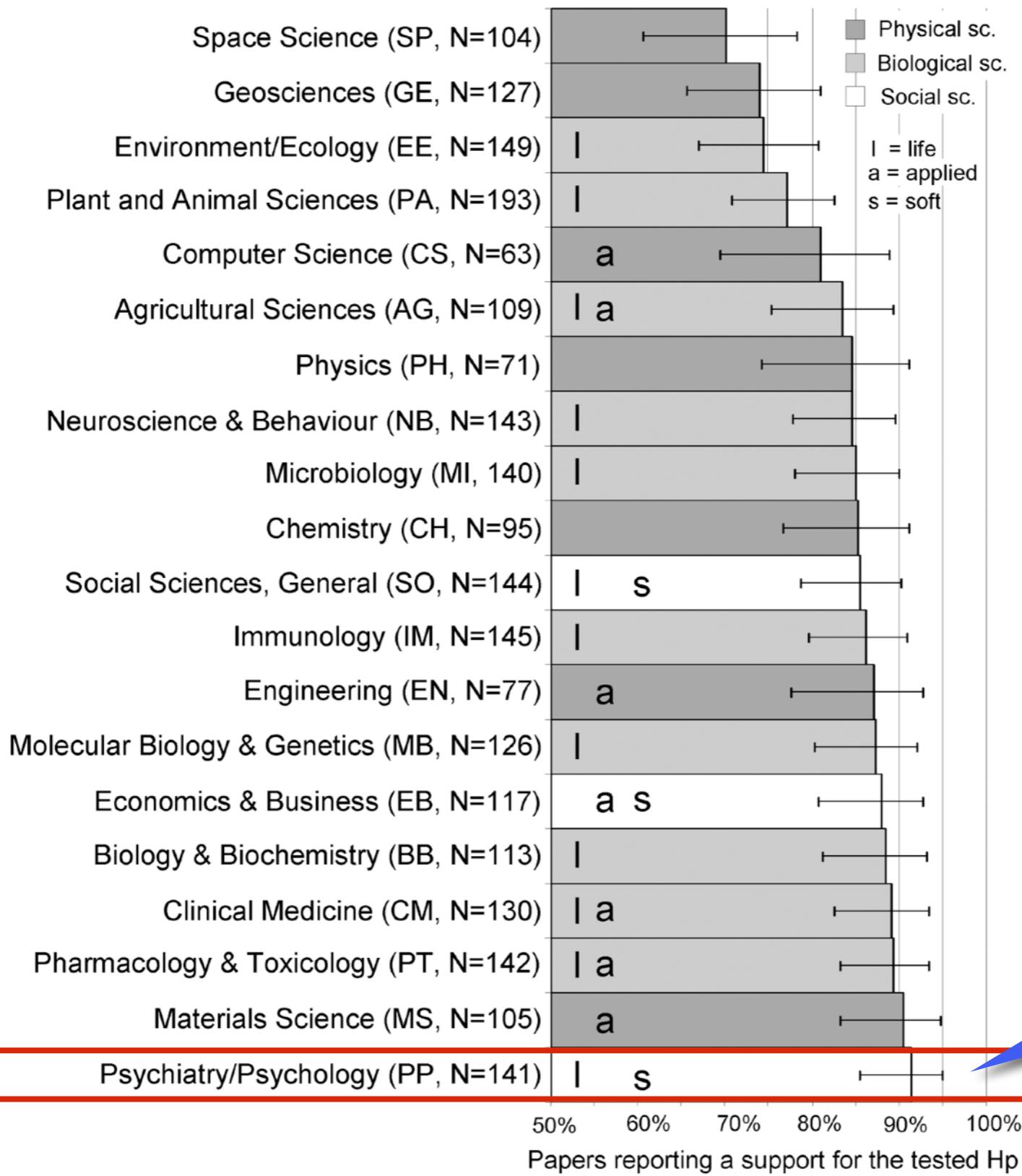
20:05 - 13. Jan. 2015

# How to become a Professor?

Actual (not desired) relevance in professorship hiring committees	Rank
<b>Number</b> of peer-reviewed publications	1
Fit of research profile to the hiring department	2
Quality of research talk	3
<b>Number</b> of publications	4
<b>Volume</b> of acquired third-party funding	5
<b>Number</b> of first authorships	6
...	...

N = 1453 psychology researchers, 66% were actually members of a professorship hiring committee.

# How to get lots of publications?



92% of published  
papers have  
significant,  
positive results

***p-hack your way  
to scientific glory!***

**p-hacking** (*n.*). Tune your data analysis in a way that you achieve a significant *p*-value in situations where it would have been non-significant.

**Questionable research practices (QRPs)** (*n.*). Practices of data collection and data analysis that are not outright fraud, but also not really kosher.

# Tool I: Outcome switching



PROJECT RESULTS TEAM BLOG FAQ

## Tracking switched outcomes in clinical trials

Here's what we found.

<b>67</b>	<b>9</b>	<b>300</b>	<b>357</b>
TRIALS CHECKED	TRIALS WERE PERFECT	OUTCOMES NOT REPORTED	NEW OUTCOMES SILENTLY ADDED

On average, each trial reported just 62.1% of its specified outcomes. On average, each trial silently added 5.3 new outcomes.

For ██████████, “the authors conducted two additional money priming studies that showed no effects, the details of which were shared with us.” and “reported nine dependent measures that were statistically affected by the manipulation in the predicted direction (one in each experiment) but did not report 19 additional measures that were statistically unchanged”.

# Tool I: Outcome switching

- 2 outcome variables:

false positive rate **5% → 9.5%**

- 5 outcome variables with one-sided testing:

false positive rate **5% → 41%**

- How prevalent is it?

- John, Loewenstein and Prelec (2012):  
66% of researchers admit having done this.

# Tool 2: Many conditions, report only those that worked

- Assess more than two conditions (and leave out conditions that are not significantly different).
- E.g., testing “high”, “medium” and “low” conditions and reporting only the results of a “high” versus “medium” comparison.
- Gives you more than one chance to find an effect. Can increases the false positive rate to **12.6%**.
- How prevalent is it?
  - 27% of researchers admit having done this (John et al., 2012).

# Tool 2: Many conditions, report only those that worked

**Best-practice example:**  
**Transforming a boring dissertation into a groundbreaking publication**



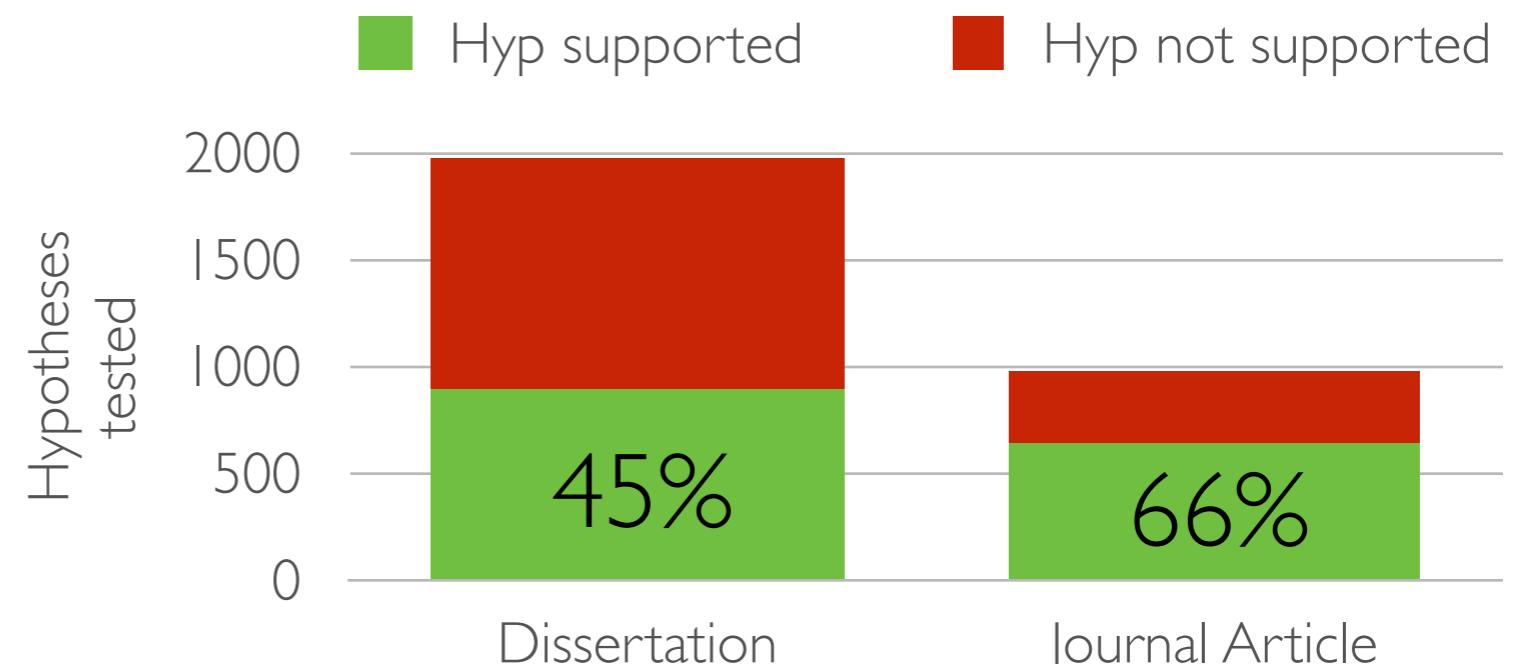
The image shows two tweets from Joe Hilgard (@JoeHilgard) on Twitter. The first tweet, posted on February 16, discusses a study involving four conditions and 415 subjects, while the manuscript version had three conditions and 140 subjects. The second tweet, also from February 16, explains that the study started with a  $2 \times 2 \times 4$  design and was simplified to a  $2 \times 3$  design.

**Joe Hilgard** @JoeHilgard · 16. Feb.  
Here's another spicy one: Thesis reports four conditions, 415 subjects. Manuscript reports three conditions, 140 subjects.

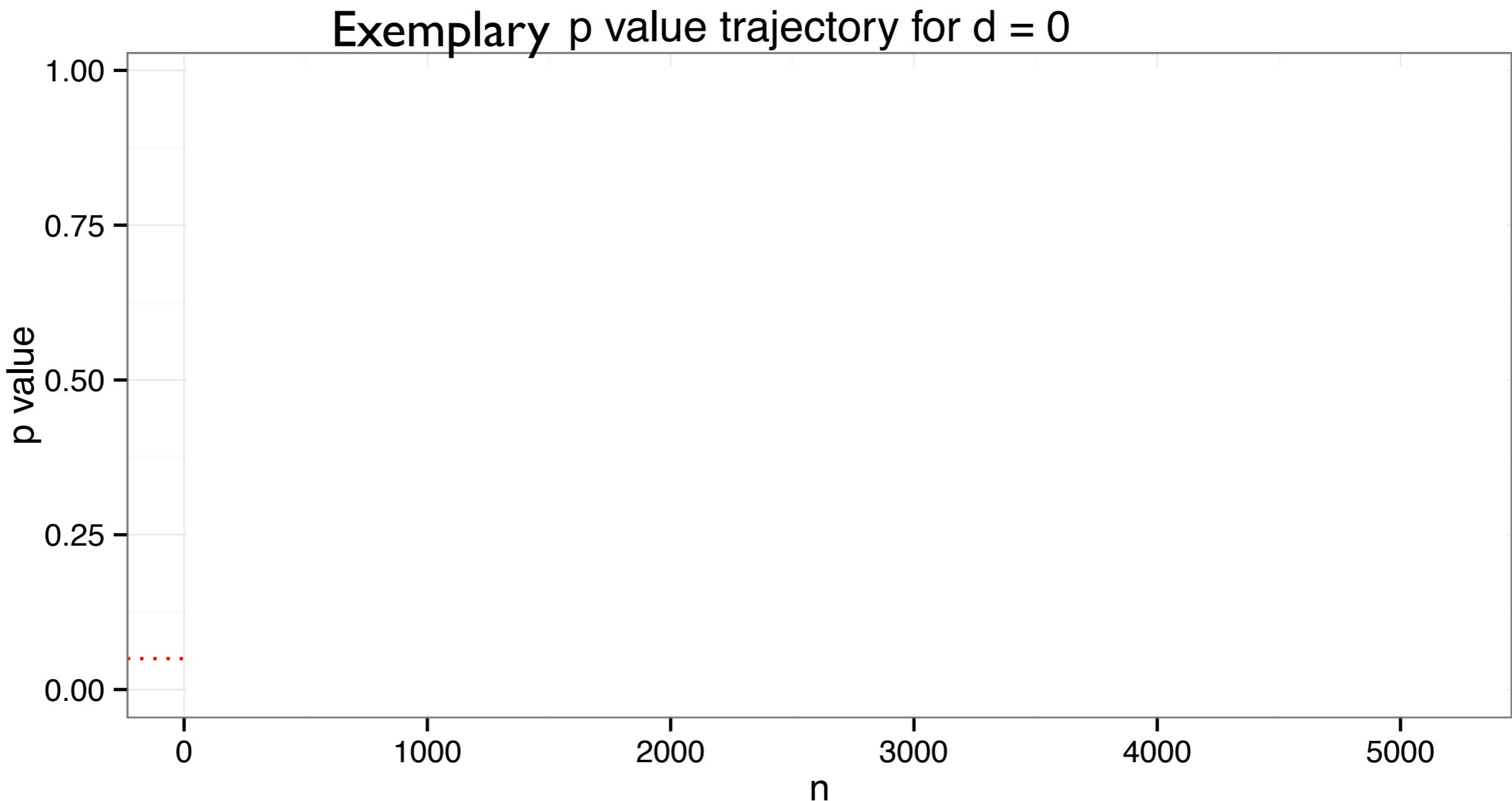
**Joe Hilgard** @JoeHilgard · 16. Feb.  
Figured it out: It started with a  $2 \times 2 \times 4$  design and worked its way down to the  $2 \times 3$  design that "worked."

## The „Chrysalis effect“ (O’Boyle et al, 2017)

- 142 dissertations that were subsequently published in a refereed journal (149 field studies and 26 experiments)

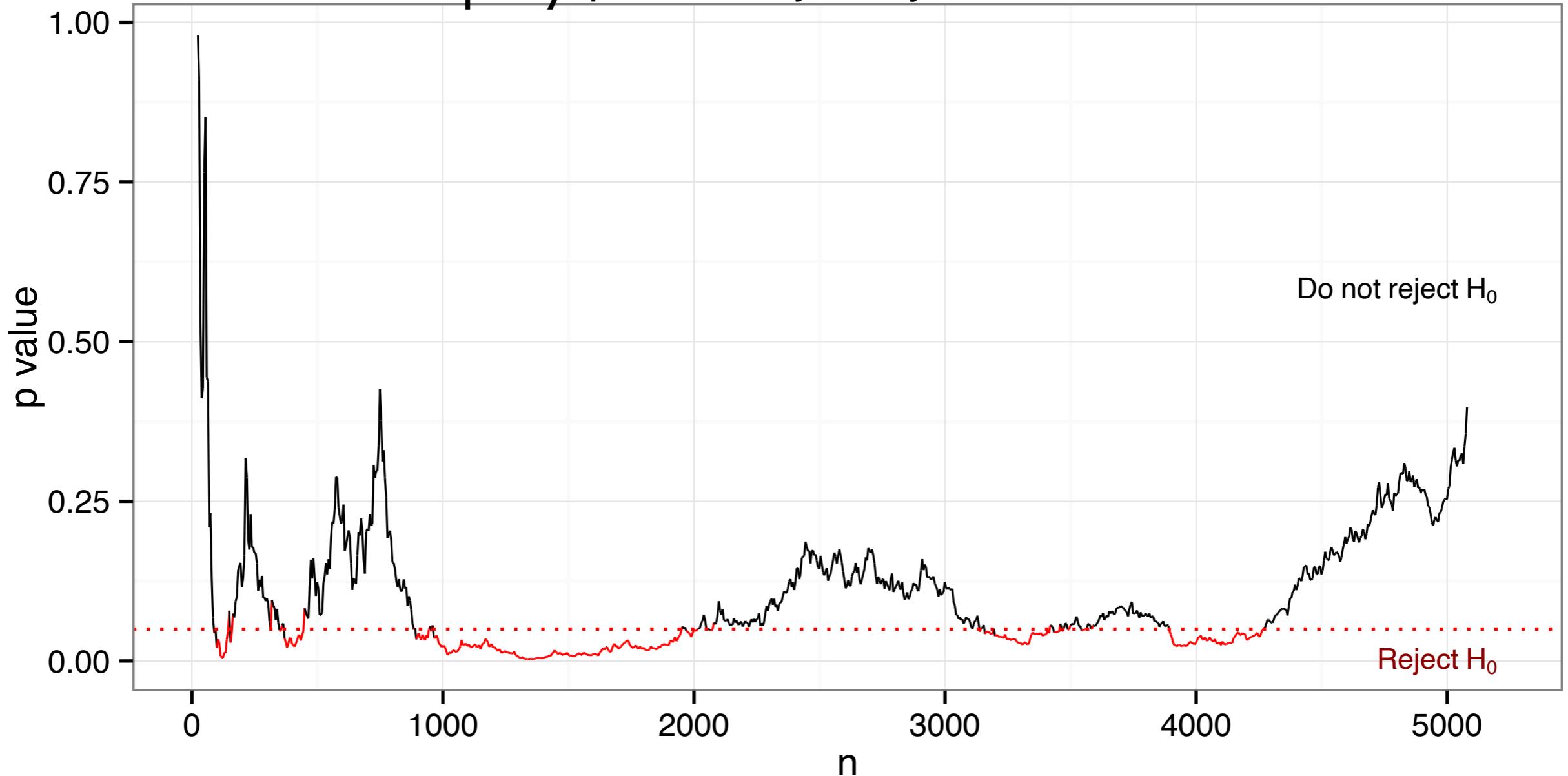


# $p$ -values under $H_0$



# Under $H_0$ , $p$ values meander infinitely

Exemplary  $p$  value trajectory for  $d = 0$



# Repeated Significance Tests on Accumulating Data

By P. ARMITAGE, C. K. MCPHERSON and B. C. ROWE

*Department of Medical Statistics and Epidemiology,  
London School of Hygiene and Tropical Medicine*

TABLE 2

*The probability of being absorbed at or before the nth observation in sampling from a normal distribution with known variance, with repeated tests at a nominal two-sided significance level  $2\alpha$  (i.e. standardized normal deviate k)†*

$2\alpha$ $k$	0·10		0·05		0·02		0·01	
	1·645		1·960		2·326		2·576	
$n$	$Q$	$S$	$Q$	$S$	$Q$	$S$	$Q$	$S$
1	0·10000	0·0970	0·05000	0·0545	0·02000	0·0230	0·01000	0·0135
2	0·16015	0·1650	0·08312	0·0885	0·0345			
3	0·20207	0·1980	0·10726	0·1115	0·0456			
4	0·23399	0·2295	0·12617	0·1260	0·0545			
5	0·25963	0·2590	0·14169	0·1420	0·0620			
160	0·63315		0·40829		0·2083			
180	0·64301		0·41677		0·2135			
200	0·65165		0·42429		0·2182			
250	0·670		0·440		0·228			
500	0·720		0·487		0·259			
750	0·746		0·513		0·276			
1,000	0·763		0·529		0·288		0·172	

With long enough sampling and optional stopping, it is guaranteed to get a significant result!

100%

# Tool 3: Optional stopping

- Collect an initial sample, analyze the results, add additional participants if not significant, stop when significance is found
- Increase twice:  $\alpha = \underline{11\%}$
- But with enough looks can be pushed to **100%**!
- How prevalent is it?
  - 70% of researchers admit having continued or stopped data collection based on looking at the interim results (John et al., 2012).

# Tool 4: Multiple comparisons in ANOVA

- ANOVA, 3 factors, full model
  - 3 main effects, 3 two-way interactions, 1 three-way interaction
- Type I error rate for at least 1 significant term?
- Well-Known: Corrections for post-hoc comparisons of levels within one factor
- Less-known: The need for correcting multiple interactions.



# Tool 5: Subgroup analyses

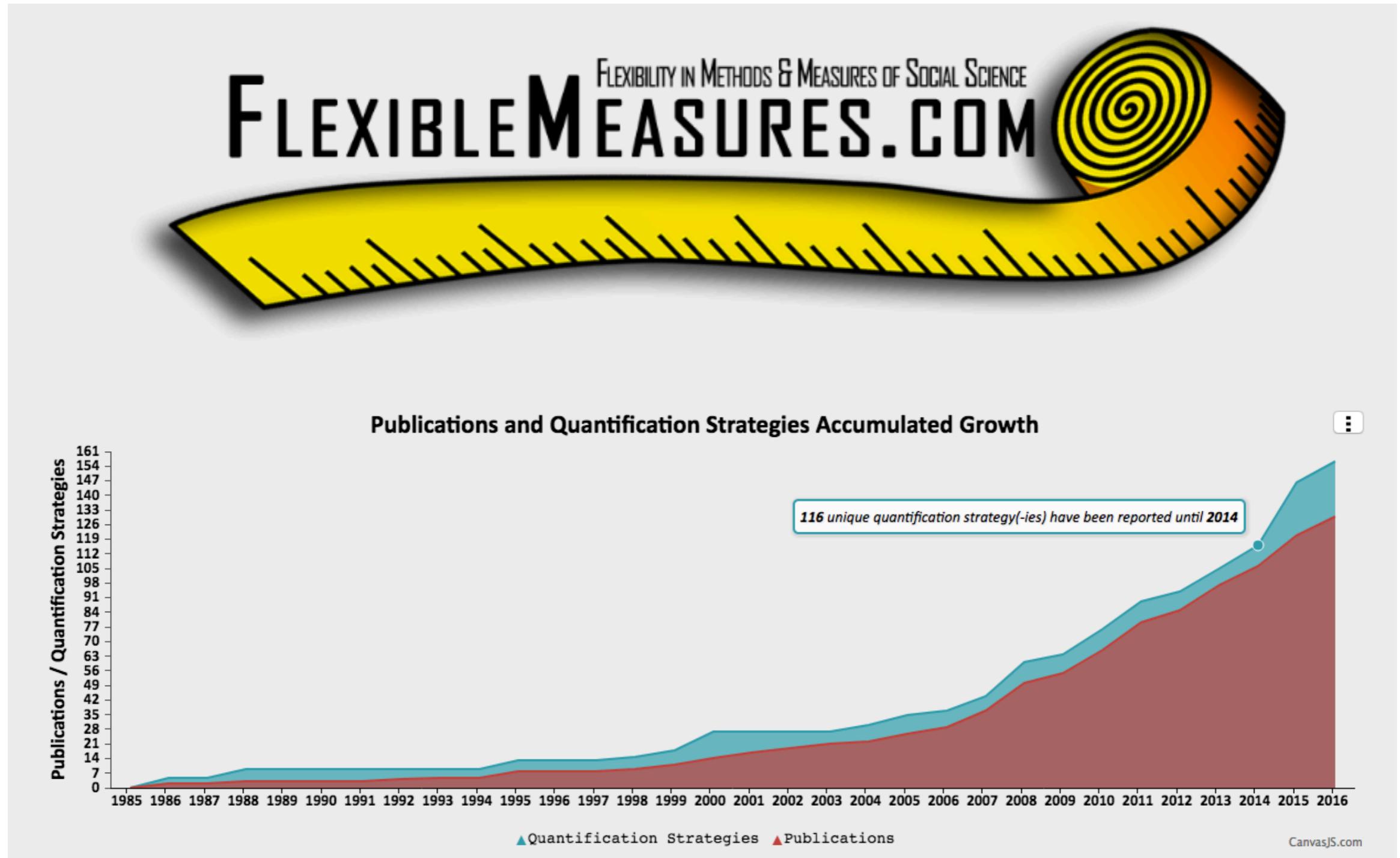
Research question: Do aggressive primes trigger aggressive behavior?

A second study in Turner, Layton, and Simons (1975) collects a larger sample of men and women driving vehicles of all years. **The design was a 2 (Rifle: present, absent) × 2 (Bumper Sticker: "Vengeance", absent) design with 200 subjects.**

→ presumably, no effect ... (yet! Do not give up so easily)

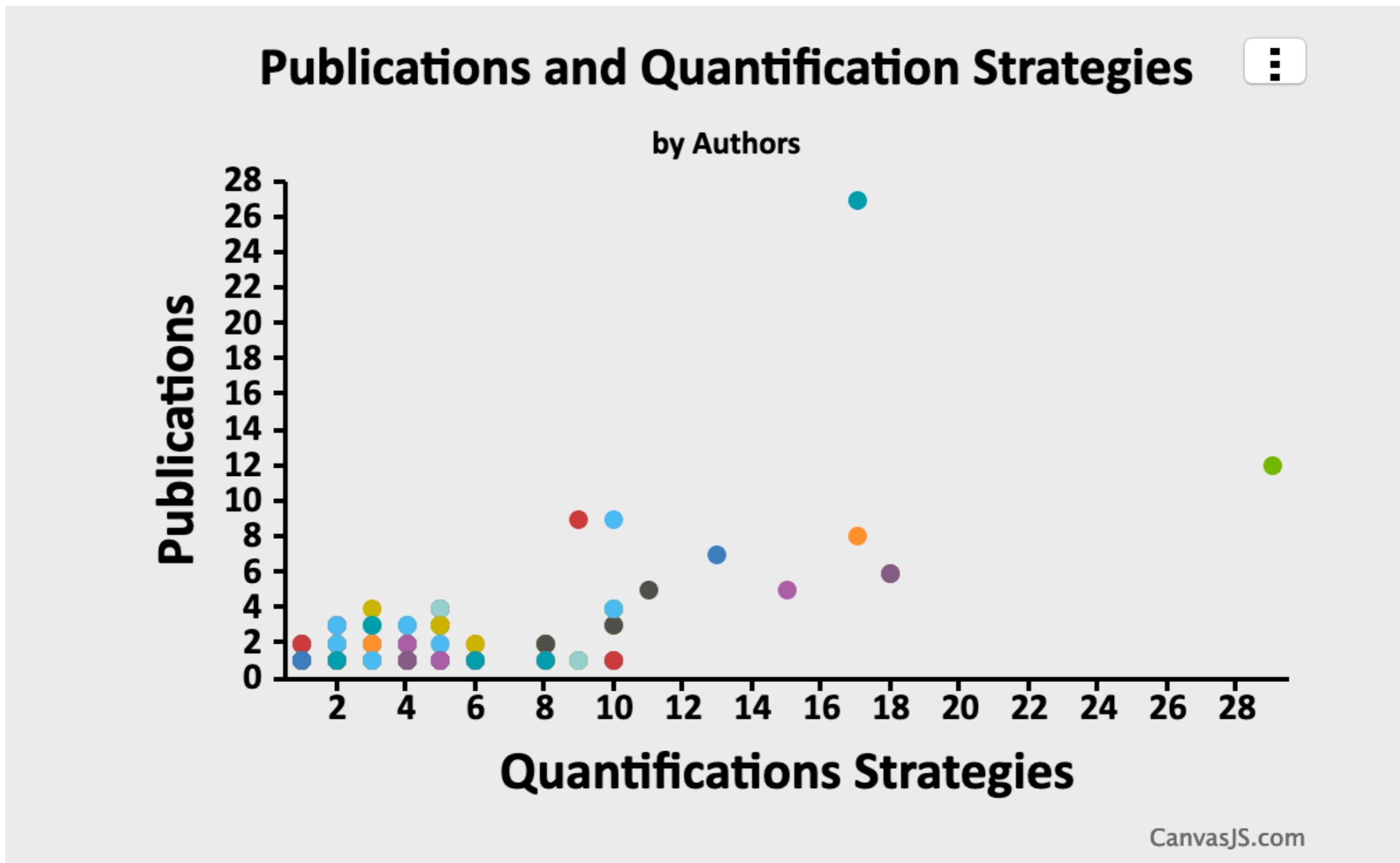
They **divide this further by driver's sex** and by a **median split on vehicle year**. They find that the Rifle/Vengeance condition increased honking relative to the other three, but only among newer-vehicle male drivers,  $F(1, 129) = 4.03, p = .047$ . But then they report that the Rifle/Vengeance condition decreased honking among older-vehicle male drivers,  $F(1, 129) = 5.23, p = .024$ ! No results were found among female drivers.

# Tool 6: Flexible measures



<http://www.flexiblemeasures.com/> by Malte Elson

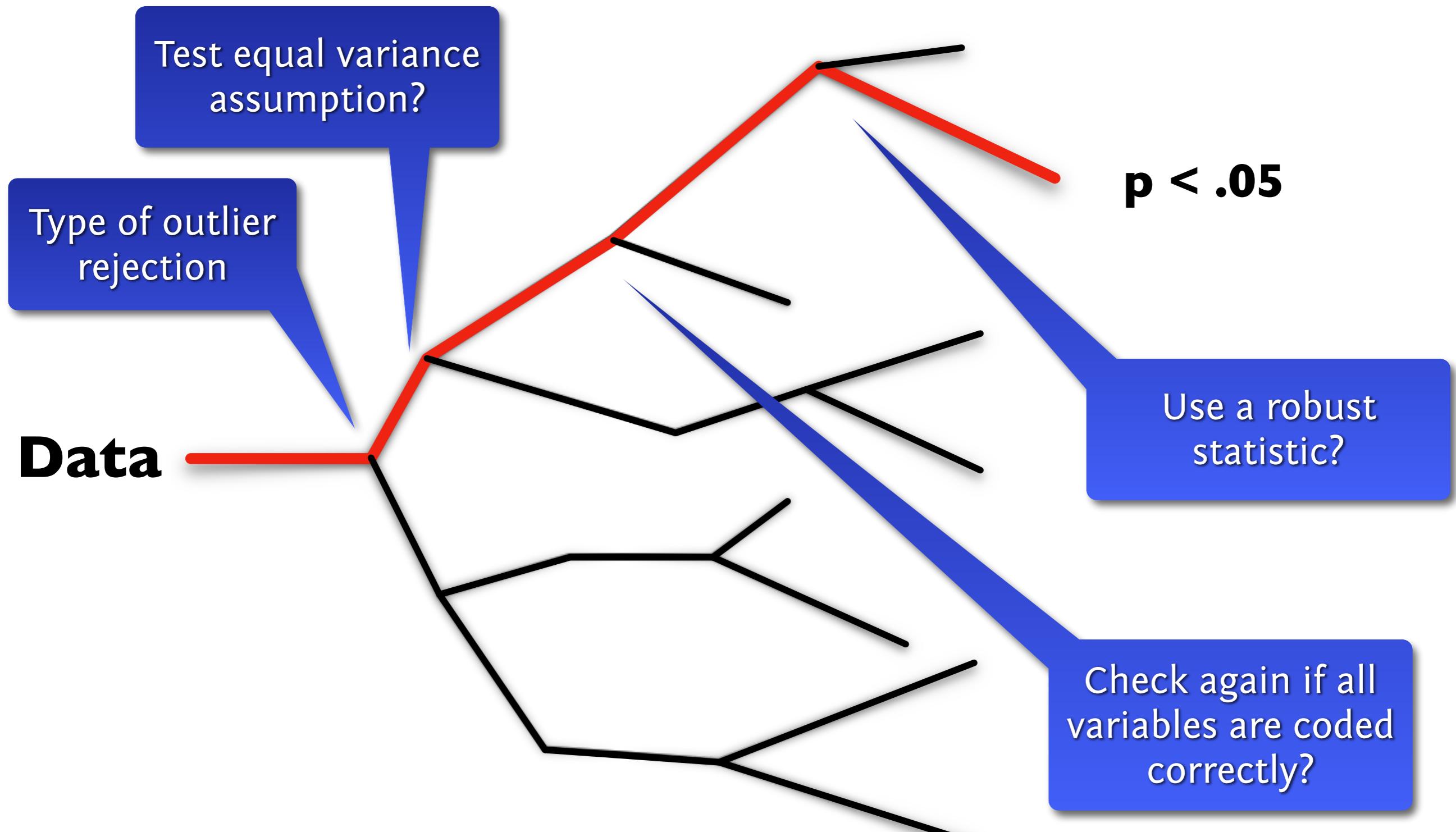
# Tool 6: Flexible measures



<http://www.flexiblemeasures.com/> by Malte Elson

# Tool 7: Explore the garden of forking paths

Andrew Gelman & Eric Loken, 2013



# Tool 8: Build the $p$ -hacking into the software!

**PNAS**

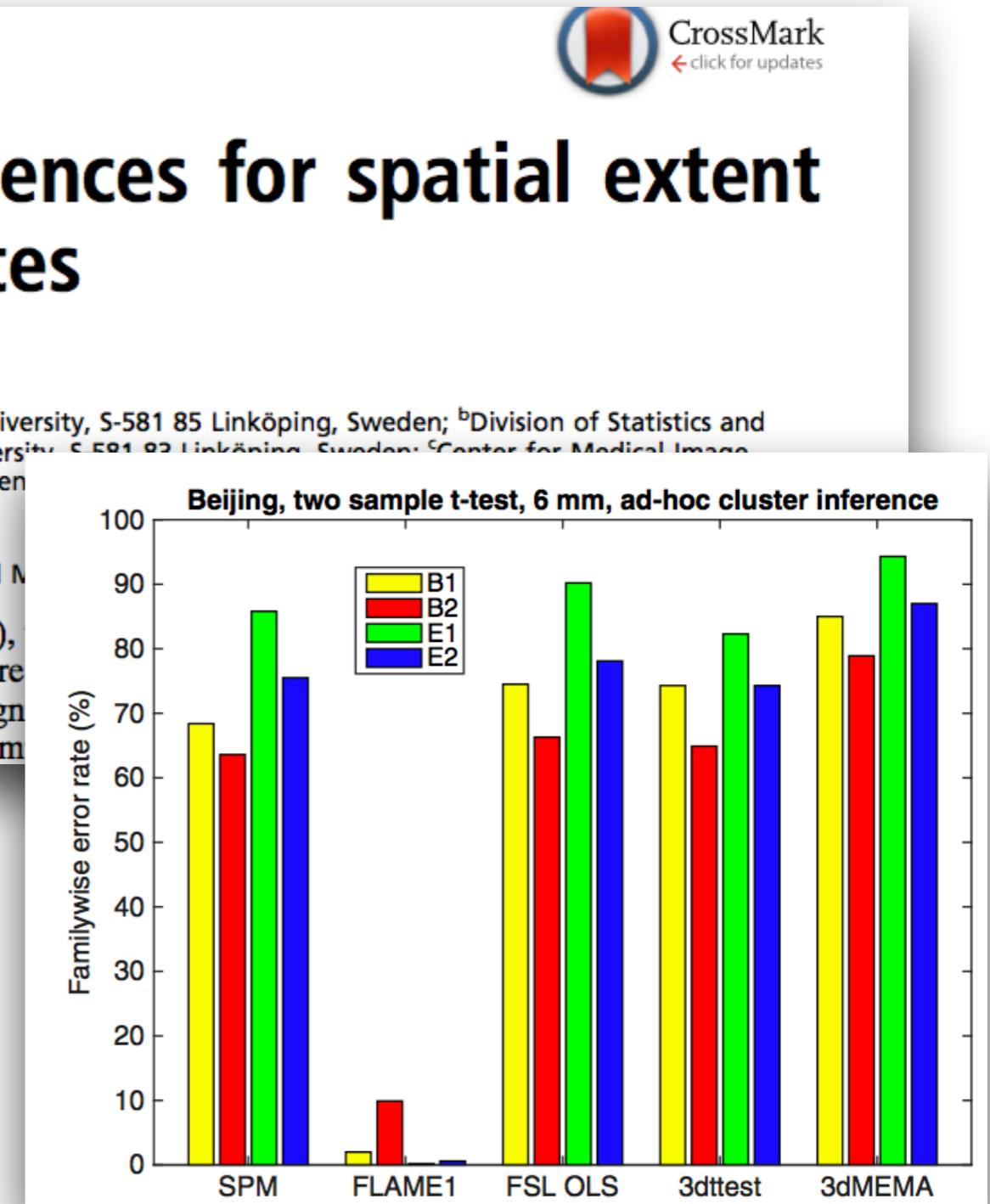
## Cluster failure: Why fMRI inferences for spatial extent have inflated false-positive rates

Anders Eklund<sup>a,b,c,1</sup>, Thomas E. Nichols<sup>d,e</sup>, and Hans Knutsson<sup>a,c</sup>

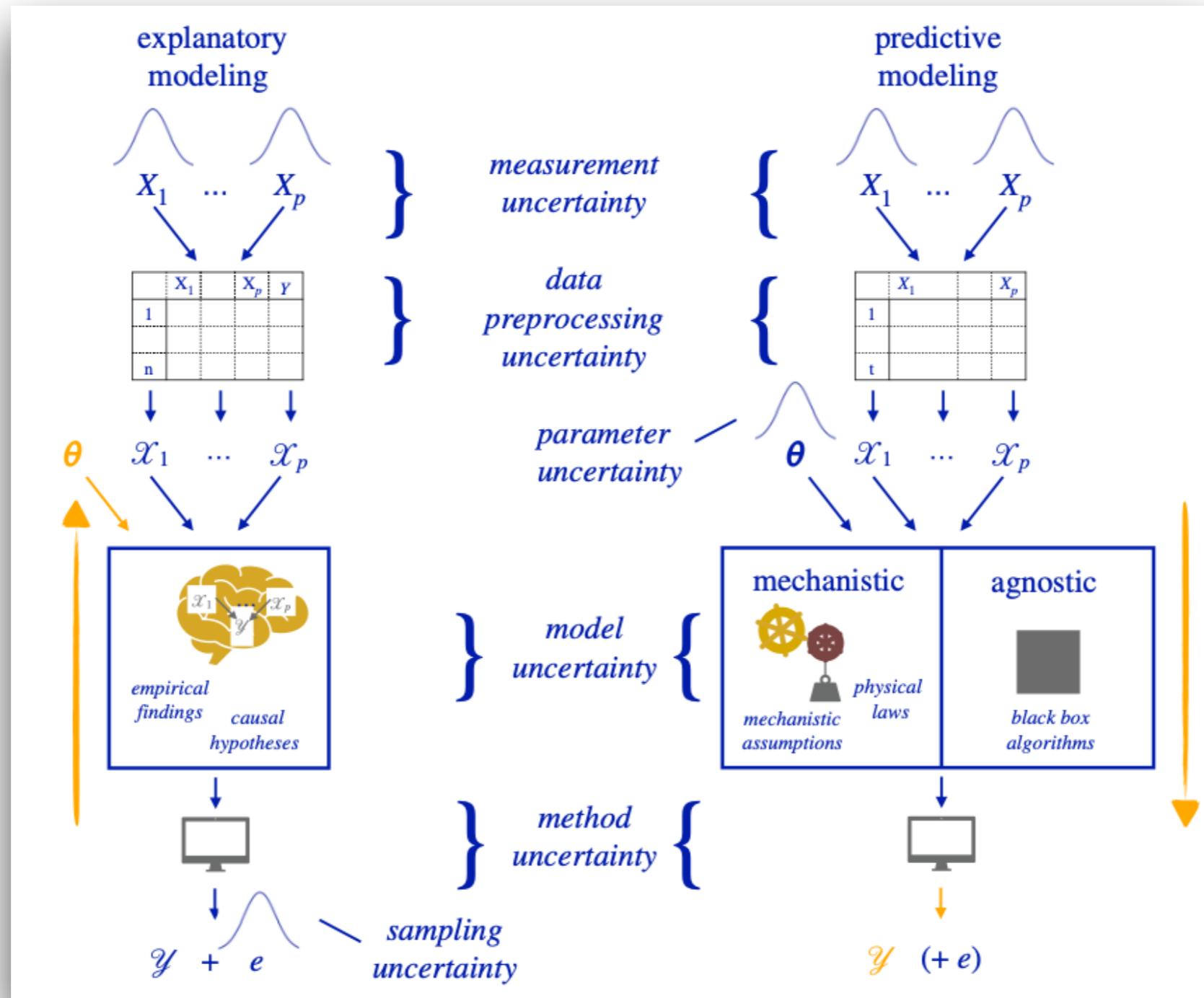
<sup>a</sup>Division of Medical Informatics, Department of Biomedical Engineering, Linköping University, S-581 85 Linköping, Sweden; <sup>b</sup>Division of Statistics and Machine Learning, Department of Computer and Information Science, Linköping University, S-581 87 Linköping, Sweden; <sup>c</sup>Center for Medical Image Science and Visualization, Linköping University, S-581 83 Linköping, Sweden; <sup>d</sup>Department of Mathematics, University College London, London WC1E 6BT, United Kingdom; and <sup>e</sup>WMG, University of Warwick, Coventry CV4 7AL, United Kingdom

Edited by Emery N. Brown, Massachusetts General Hospital, Boston, MA, and approved November 13, 2013.

The most widely used task functional magnetic resonance imaging (fMRI) analyses use parametric statistical methods that depend on a variety of assumptions. In this work, we use real resting-state data and a total of 3 million random task group analyses to compute



# The multiplicity of analysis strategies jeopardizes replicability: Lessons learned across disciplines



*Freedom is nothing else but a chance to be better.*

*Albert Camus (1913 - 1960)*

A black and white photograph of Albert Camus, a French-Algerian writer, philosopher, and anti-colonial activist. He is shown from the chest up, wearing a dark jacket over a light-colored shirt. His hair is neatly combed, and he has a thoughtful expression, looking slightly to his left. The background is blurred, showing what appears to be an interior room.

Ok, let's celebrate  
some researcher degrees  
of freedom to be better  
at research!

# Train your skills!

p-hacker: Train your p-hacking skills!

Manual

New study Now: p-hack!

Settings for initial data collection:

Name for experimental group  
Elderly priming

Name for control group  
Control priming

Initial # of participants in each group  
20

True effect in population  
0

Number of DVs  
4

Run new experiment  
(Discards previous data)

Tests for each DV

Name	N	Statistic	p-Value	Sign.	Actions
DV1	40	$F(1, 38) = 1.02$	$p = .318$	ns	Save
DV2	40	$F(1, 38) = 1.32$	$p = .257$	ns	Save
DV3	40	$F(1, 38) = 1.37$	$p = .249$	ns	Save
DV4	40	$F(1, 38) = 1.24$	$p = .272$	ns	Save
DV_all	39	$F(1, 37) = 3.79$	$p = .059$	ns	Save

Choose DV to plot  
DV\_all

http://shinyapps.org/apps/p-hacker/

The impact of  $p$ -hacking on the  
rate of significant results

# It is done . . .

**Table 1.** Biostatistician-Reported Frequency and Severity Rating of Requests for Inappropriate Analysis and Reporting ( $n = 390$ )\*

Violation Request	Respondents Rating the Item as "Most Severe," %†	Reported Requests During the Past 5 Years, %		
		0	1-9	≥10
Falsify the statistical significance (such as the $P$ value) to support a desired result	84	97	2	1
Change data to achieve the desired outcome (such as the prevalence rate of cancer or another disease)	84	93	7	-
Remove or alter some data records (observations) to better support the research hypothesis	80	76	22	2
Interpret the statistical findings on the basis of expectations, not the actual results	68	70	28	2
Do not fully describe the treatment under study because protocol was not exactly followed	62	85	15	-
Do not report the presence of key missing data that could bias the results	68	76	23	1
Ignore violations of assumptions because results may change to negative	64	71	28	1
Modify a measurement scale to achieve some desired results rather than adhering to the original scale as validated	55	79	20	1
Report power on the basis of a post hoc calculation, but make it seem like an <i>a priori</i> statement	54	76	23	2
Request to not properly adjust for multiple testing when "a priori, originally planned secondary outcomes" are shifted to an "a posteriori primary outcome status"	56	80	18	2
Conduct too many post hoc tests, but purposefully do not adjust $\alpha$ levels to make results look more impressive than they really are	54	60	36	4
Remove categories of a variable to report more favorable results	48	68	31	1
Do not mention interim analyses to avoid "too much testing"	50	81	18	1
Report results before data have been cleaned and validated	48	56	39	5
Do not discuss the duration of follow-up because it was inconsistent	45	84	15	1
Stress only the significant findings, but underreport nonsignificant ones	42	45	48	7
Do not report the model statistics (including effect size in ANOVA or $R^2$ in linear regression) because they seemed too small to indicate any meaningful changes	42	76	23	1
Do not show plot because it did not show as strong an effect as you had hoped	33	58	39	3

ANOVA = analysis of variance.

\* Based on findings from questions 1-18 of the Bioethical Issues in Biostatistical Consulting Questionnaire, which asked biostatisticians "to estimate the number of times—during the past 5 years—that you, personally, have been DIRECTLY asked to do this." Data are presented in decreasing order by the percentage of respondents with a perceived severity score of 4 or 5.

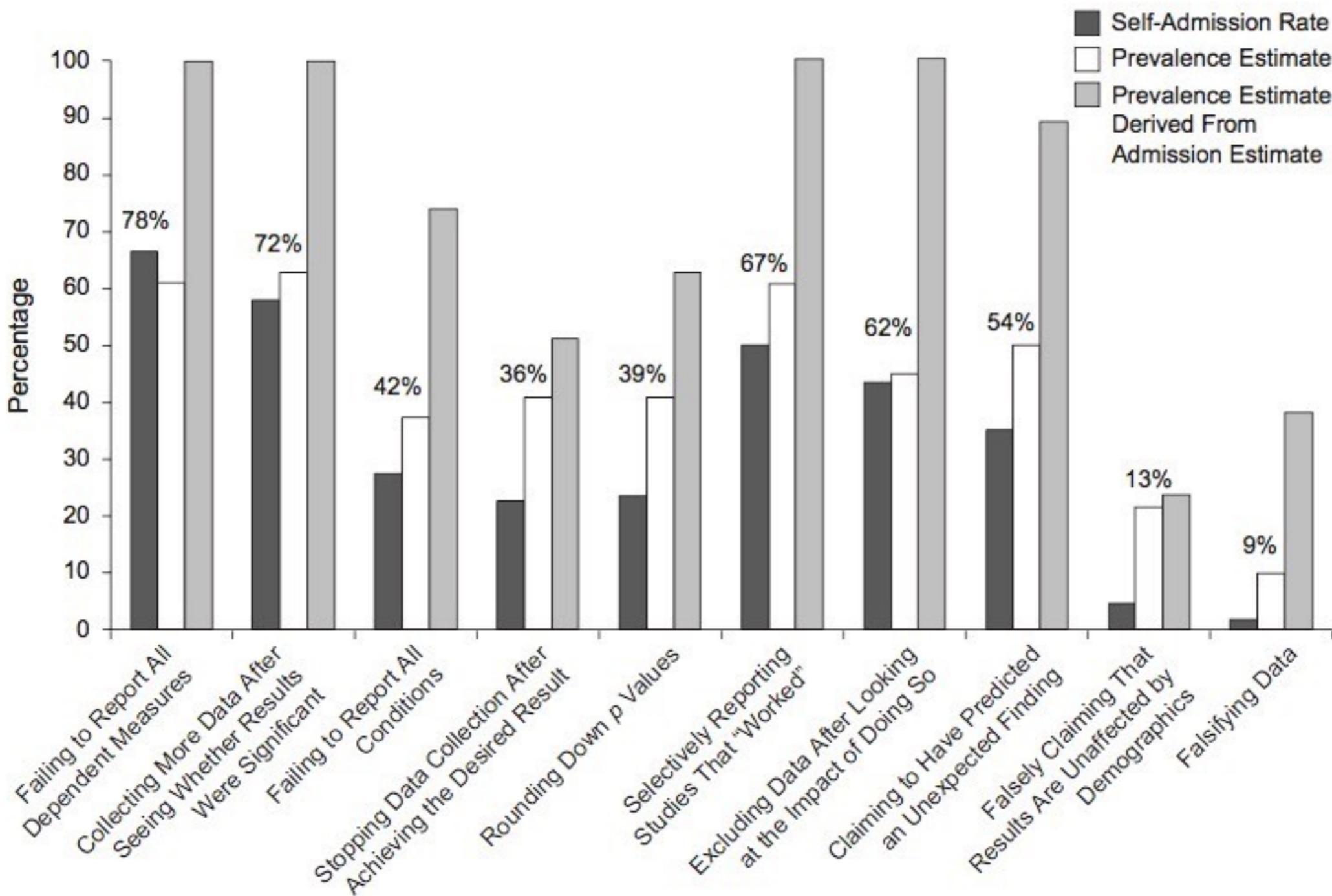
† Items were defined as "most severe" if respondents ranked the severity as 4 or 5 on a scale of 0-5.

# Measuring the Prevalence of Questionable Research Practices With Incentives for Truth Telling *in psychology*

Psychological Science  
23(5) 524–532  
© The Author(s) 2012  
Reprints and permission:  
[sagepub.com/journalsPermissions.nav](http://sagepub.com/journalsPermissions.nav)  
DOI: 10.1177/0956797611430953  
<http://pss.sagepub.com>  


Leslie K. John<sup>1</sup>, George Loewenstein<sup>2</sup>, and Drazen Prelec<sup>3</sup>

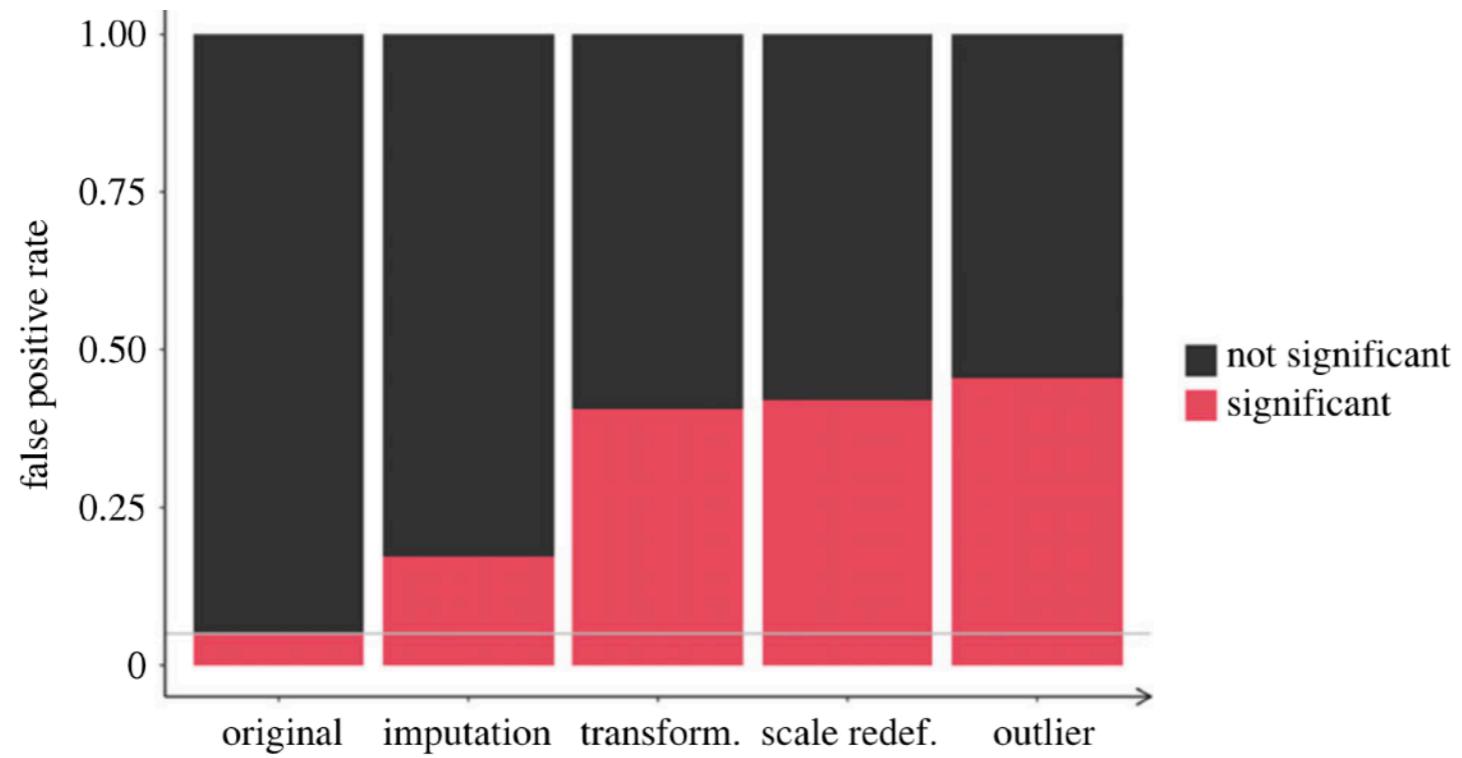
<sup>1</sup>Marketing Unit, Harvard Business School; <sup>2</sup>Department of Social & Decision Sciences, Carnegie Mellon University; and <sup>3</sup>Sloan School of Management and Departments of Economics and Brain & Cognitive Sciences, Massachusetts Institute of Technology



# Survey among 6,813 academic researchers in The Netherlands: Self-reported prevalence of fabrication and falsification in the last 3 years

			Disciplinary field		
QRP	Description (In the last three years.)	Life and medical sciences	Social and behavioural sciences	Natural and engineering sciences	Arts and humanities
<b>Fabrication</b>	Making up of data or results				
<b>Falsification</b>	Manipulating research materials, data or results				
<b>Any FF</b>	Fabrication and/or Falsification				

# How bad can it be?



- Doing some of these *questionable research practices* (QRPs) in combination can raise false positive rate from 5% to **> 50%**!
- QRPs corrupt the logic of the  $p$ -value and “renders the reported  $p$ -values essentially uninterpretable.”



# P-hacking under $H_1$

- From a statistical point of view,  $p$ -hacking increases your statistical power

$$Pr(p < .05 | H_1, phack) > Pr(p < .05 | H_1)$$

- For example:
  - Meta-analysis with  $k = 10$  studies, true effect is  $\delta = 0.2$ , typical sample sizes
  - Power **without**  $p$ -hacking in primary studies: **53%**
  - Power **with**  $p$ -hacking in primary studies: **76%**

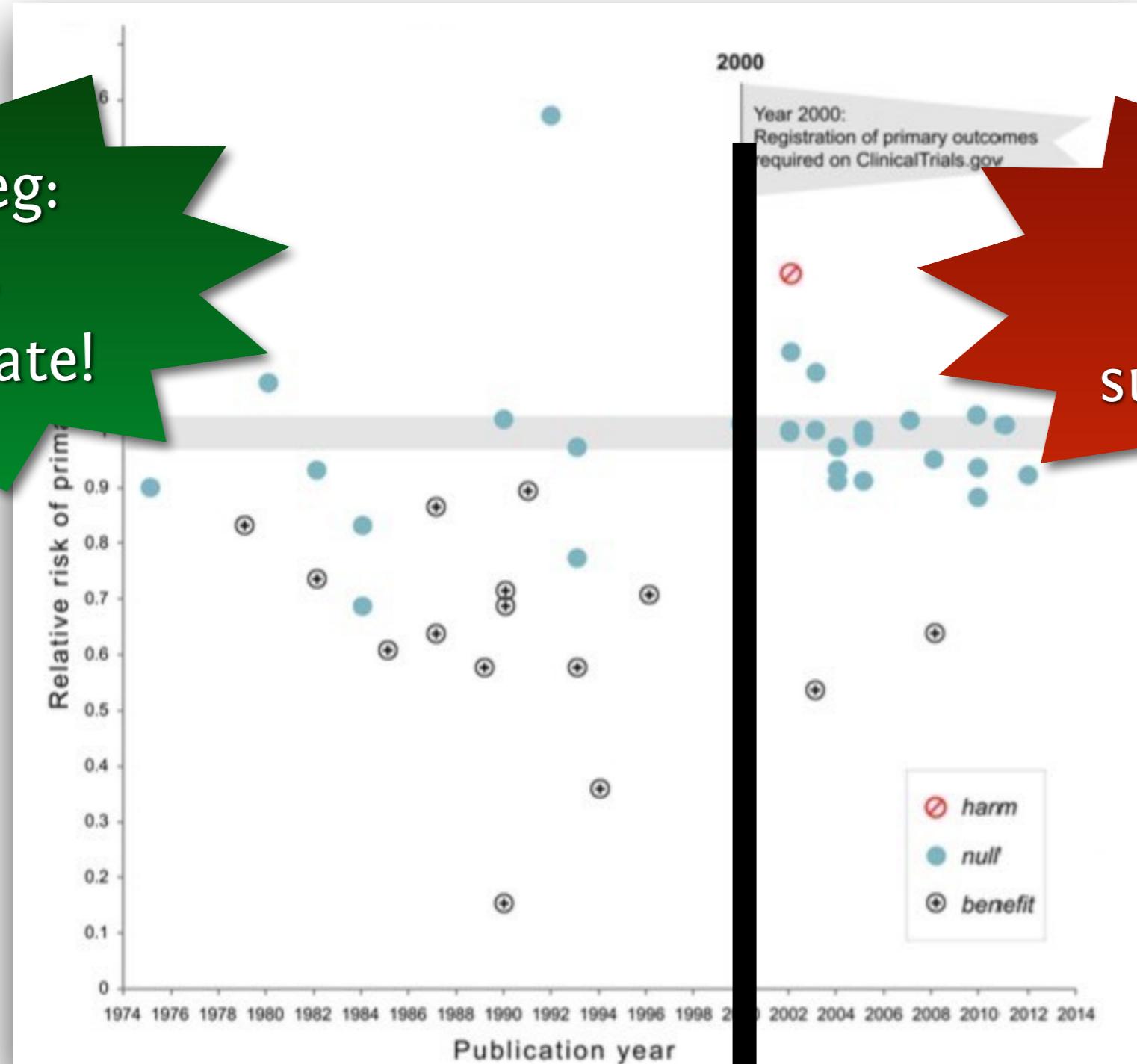
# Things to avoid

# Anti-tool: Pre-registration stops $p$ -hacking



no prereg:  
**57%**  
success rate!

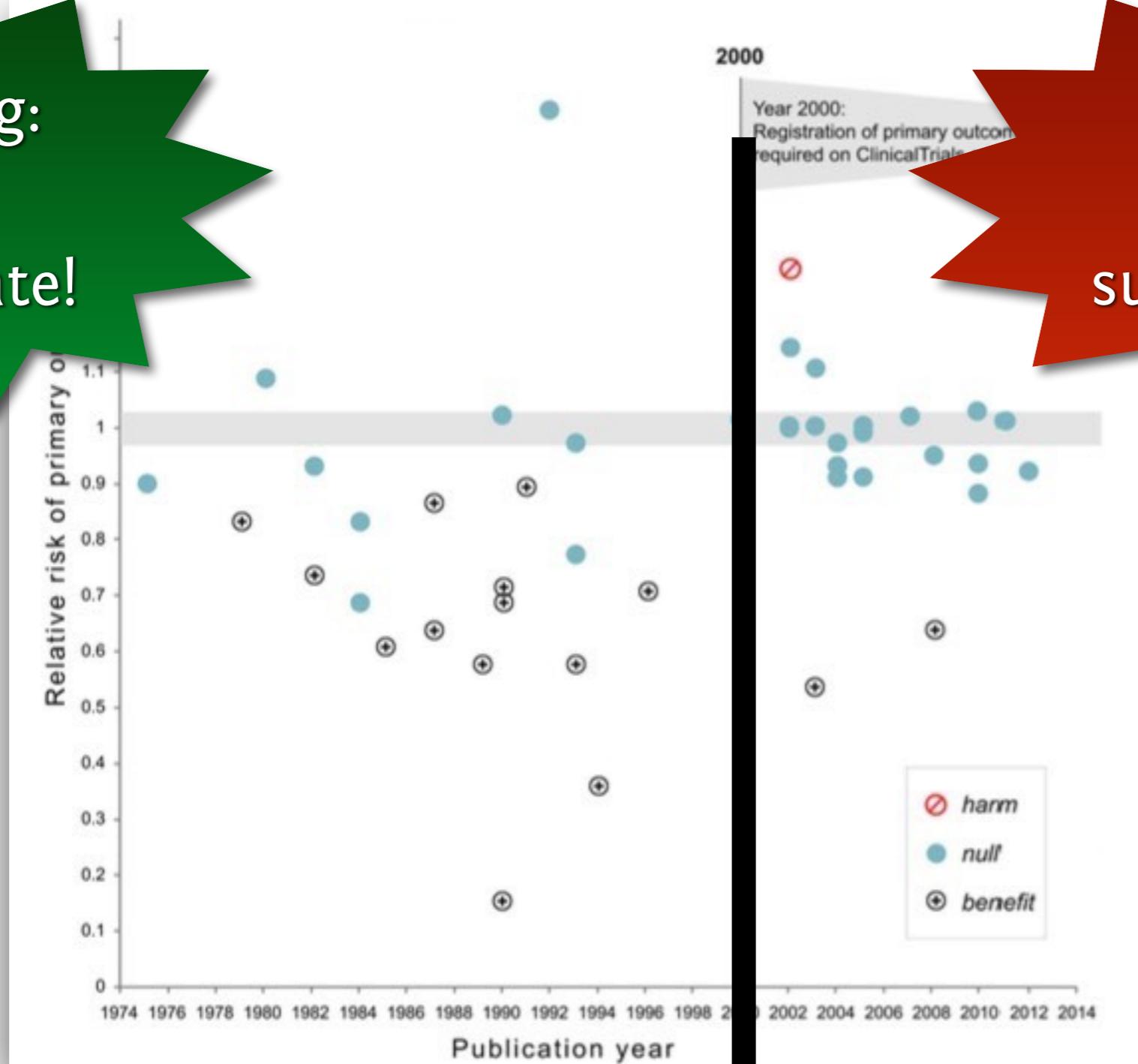
prereg:  
**8%**  
success rate...



# Tool 9: Do **not** pre-register!

no prereg:  
57%  
success rate!

prereg:  
8%  
success rate...



# Tool 10: Do **not** share open data

## **Revisiting the Power Pose Effect: How Robust Are the Results Reported by Carney, Cuddy, and Yap (2010) to Data Analytic Decisions?**

Marcus Credé<sup>1</sup> and Leigh A. Phillips<sup>1</sup>

Social Psychological and Personality Science  
1-7  
© The Author(s) 2017  
Reprints and permission:  
[sagepub.com/journalsPermissions.nav](http://sagepub.com/journalsPermissions.nav)  
DOI: 10.1177/1948550617714584  
[journals.sagepub.com/home/spp](http://journals.sagepub.com/home/spp)  


- A “multiverse analysis” (Steegen, Tuerlinchx, Gelman, & Vanpaemel, 2016): Report results for all plausible analytical decisions
- Check robustness of results: Do several analytical paths lead to comparable conclusions?
- Based on open data by Carney et al. (2010)

**Table I.** Multiverse Analysis for the Effect of Power Posing on Testosterone.

Gender Effect	Control Variables	Outlier Identification: Entire Sample (N = 39)		Outlier Identification: Test. Conditioned on Gender (N = 41)		Outlier Identification: Multivariate or No Exclusion (N = 42)	
		DV: T2 Test.	DV: Δ in Test.	DV: T2 Test.	DV: Δ in Test.	DV: T2 Test.	DV: Δ in Test.
Combined	Gender		.047 ( <i>p</i> = .19)		.019 ( <i>p</i> = .39)		.036 ( <i>p</i> = .23)
Combined	Gender and T1 test.	.029 ( <i>p</i> = .31)		.042 ( <i>p</i> = .21)		.055 ( <i>p</i> = .15)	
Combined	Gender and T1 cort.		.045 ( <i>p</i> = .21)		.017 ( <i>p</i> = .43)		.018 ( <i>p</i> = .42)
Combined	Gender, T1 test., and T1 cort.	.037 ( <i>p</i> = .26)		.040 ( <i>p</i> = .23)		.043 ( <i>p</i> = .21)	
Combined	T1 cort. and T2 cort.		.089 ( <i>p</i> = .07)		.038 ( <i>p</i> = .23)		.037 ( <i>p</i> = .24)
Combined	Gender, T1 test., T1 cort., and T2 cort.	<b>.123 (<i>p</i> = .04)</b>		.099 ( <i>p</i> = .06)		.102 ( <i>p</i> = .051)	
Men only	No controls		.192 ( <i>p</i> = .13)		.047 ( <i>p</i> = .44)		.096 ( <i>p</i> = .24)
Men only	T1 test.	.000 ( <i>p</i> = .96)		.073 ( <i>p</i> = .35)		.101 ( <i>p</i> = .25)	
Men only	T1 cort.		.184 ( <i>p</i> = .17)		.121 ( <i>p</i> = .22)		.063 ( <i>p</i> = .37)
Men only	T1 test. and T1 cort.						
Men only	T1 cort. and T2 cort.						
Men only	T1 test., T1 cort., and T2 cort.						
Women only	No controls						
Women only	T1 test.						
Women only	T1 cort.						
Women only	T1 test. and T1 cort.	.023 ( <i>p</i> = .48)		.023 ( <i>p</i> = .48)		.023 ( <i>p</i> = .48)	
Women only	T1 cort. and T2 cort.		.077 ( <i>p</i> = .19)		.077 ( <i>p</i> = .19)		.077 ( <i>p</i> = .19)
Women only	T1 test., T1 cort., and T2 cort.	.167 ( <i>p</i> = .053)		.167 ( <i>p</i> = .053)		.167 ( <i>p</i> = .053)	

Note. Entries are partial  $\eta^2$  values and (in parentheses) the associated *p* value. The entry in boldface is the effect for the analyses originally reported in the Carney, Cuddy, and Yap (2010) paper. Blank entries mean that the analyses would not be recommended for reasons described in the text. The number of women was constant across the three outlier strategies. DV = dependent variable; Test. = testosterone; cort. = cortisol; T1 = premanipulation; T2 = postmanipulation.

Of 54 plausible analyses exactly **one** was significant.  
Guess which has been reported in the original paper?

# Open Letter by Dana Carney (2016)

5. Initially, the primary DV of interest was risk-taking. We ran subjects in chunks and checked the effect along the way. It was something like 25 subjects run, then 10, then 7, then 5. Back then this did not seem like p-hacking. It seemed like saving money (assuming your effect size was big enough and p-value was the only issue).
6. Some subjects were excluded on bases such as “didn’t follow directions.” The total number of exclusions was 5. The final sample size was  $N = 42$ .
7. The cortisol and testosterone data (in saliva at that point) were sent to Salimetrics (which was in State College, PA at that time). The hormone results came back and data were analyzed.
8. For the risk-taking DV: One p-value for a Pearson chi square was .052 and for the Likelihood ratio it was .05. The smaller of the two was reported despite the Pearson being the more ubiquitously used test of significance for a
10. The self-report DV was p-hacked in that many different power questions were asked and those chosen were the ones that “worked.”



<https://www.youtube.com/watch?v=ZaNtz76dNSI&sns=em>

# (Un)Intentional?

- Intentional?
  - Evil researcher who only cares about his/her career and not at all about truth-seeking?
- Unintentional?
  - Wrong education?
  - Wrong/uncritical standards of the field?
  - Pushed by supervisors, reviewers, or editors?
    - <http://bulliedintobadscience.org/>
- My guess is that really most  $p$ -hacking occurs unintentional / in best faith.
- Distorting effects on the published record are probably comparable, but the ethical evaluations differs strongly.

# Lower the threshold?

comment

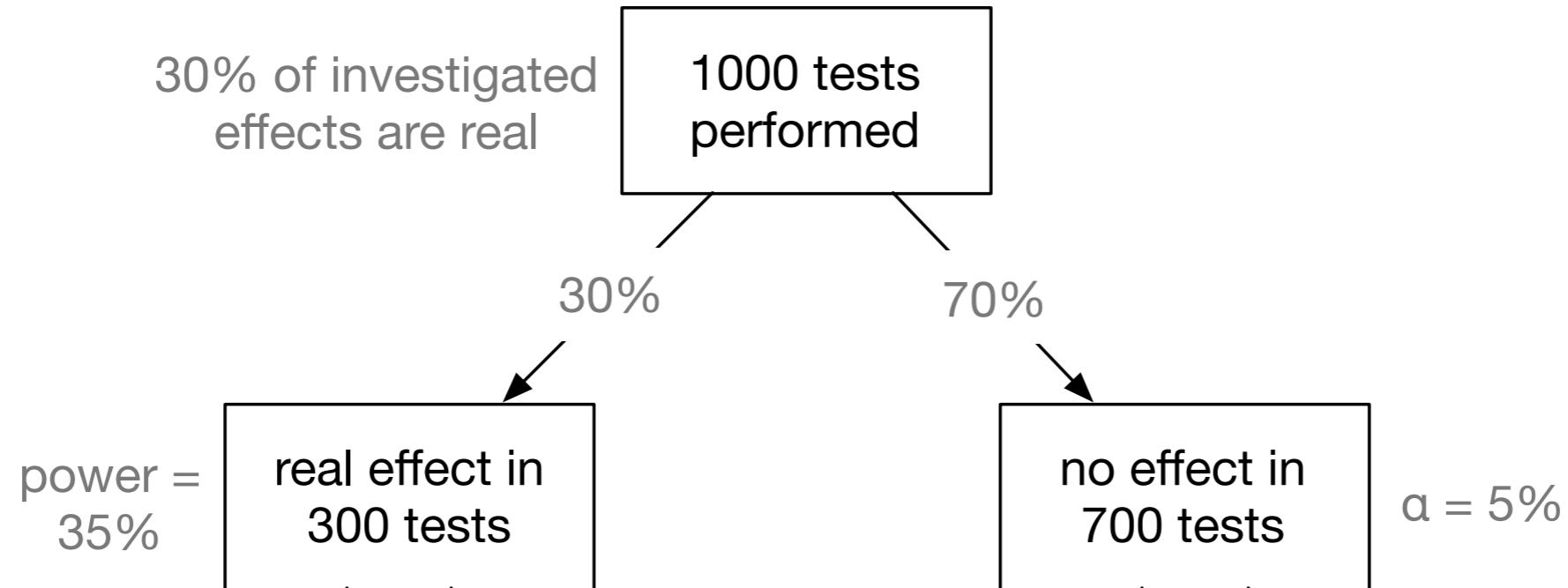
## Redefine statistical significance

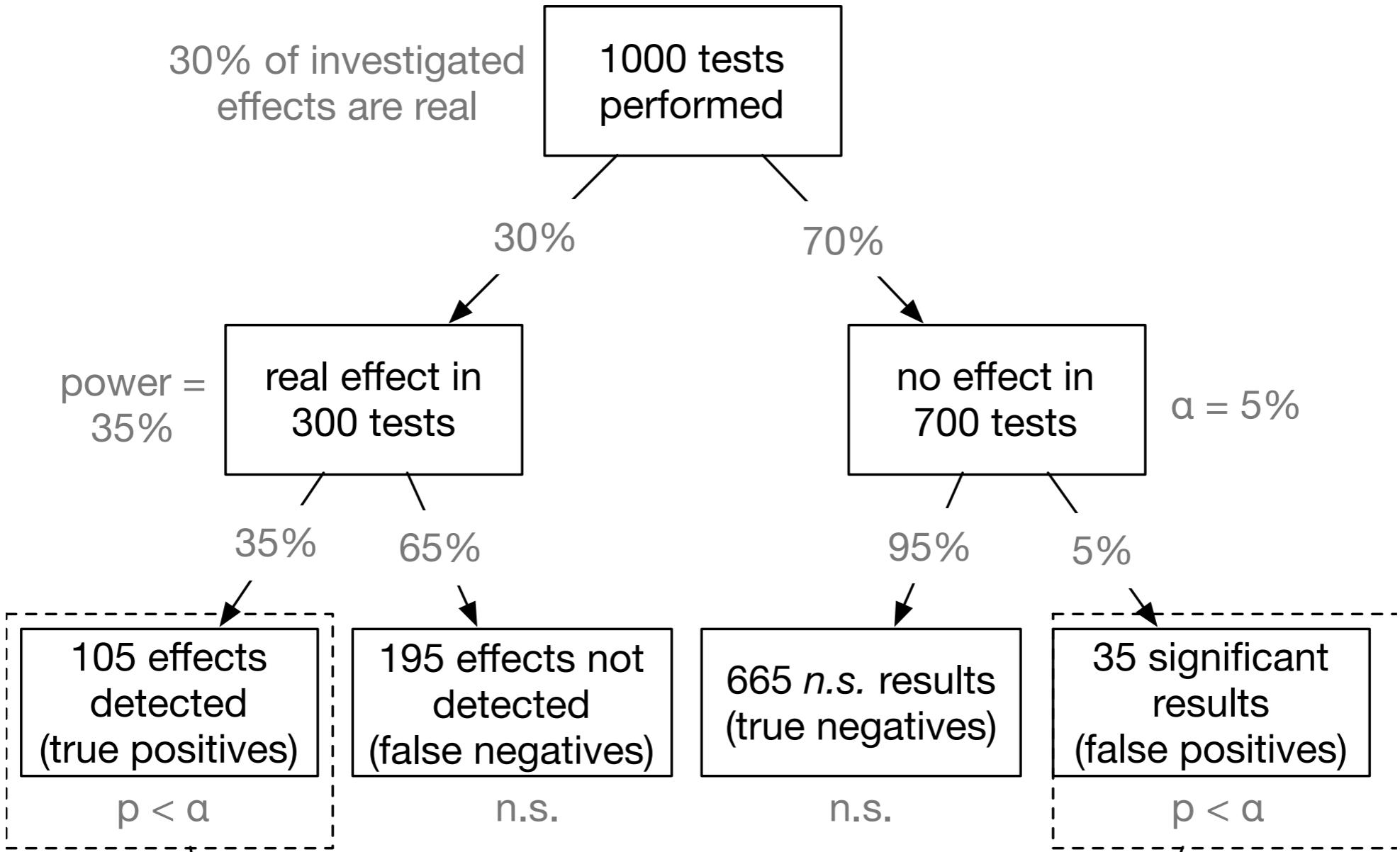
We propose to change the default  $P$ -value threshold for statistical significance from 0.05 to 0.005 for claims of new discoveries.

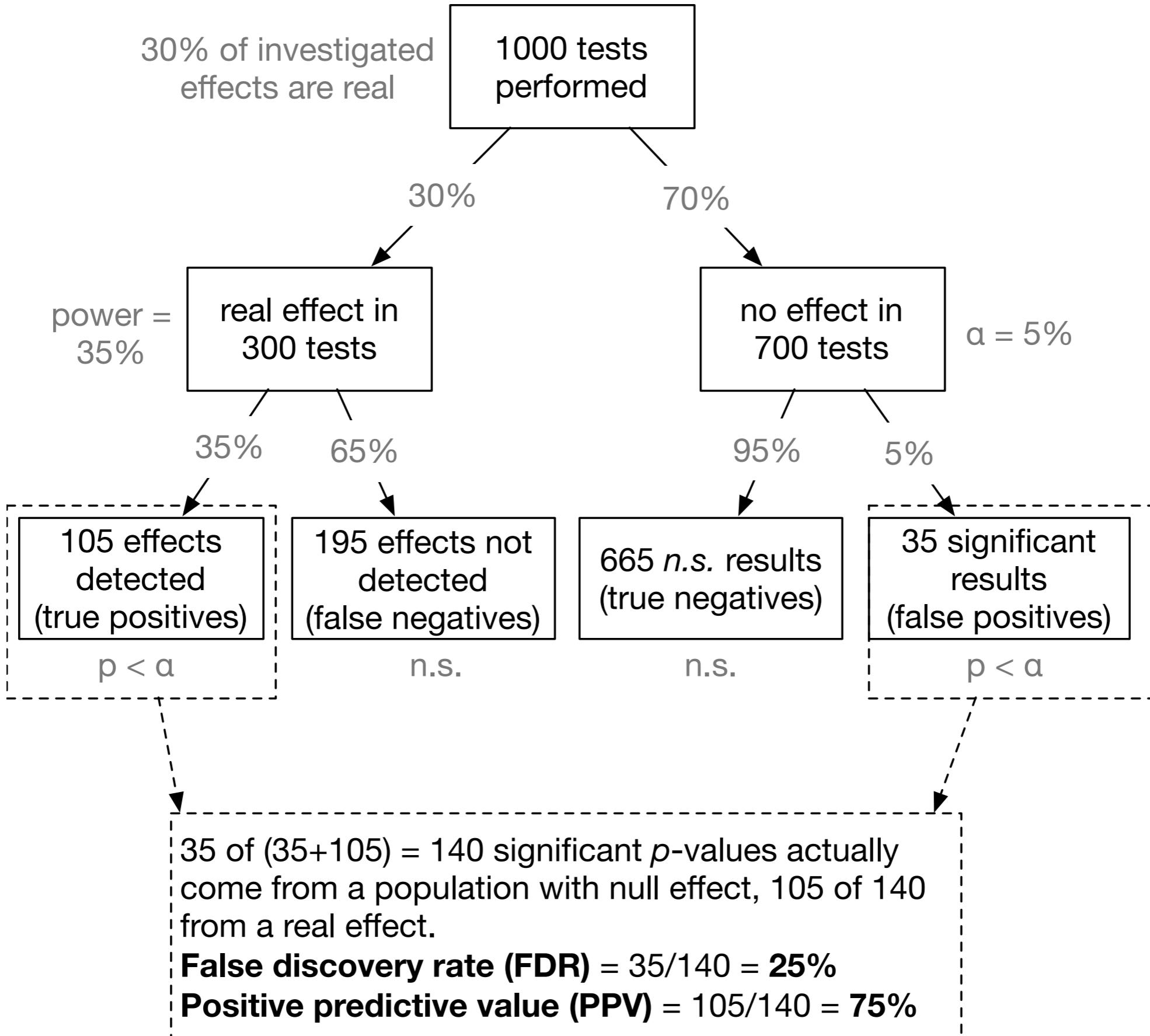
- Suggestion:
  - Lower the default threshold from .05 to .005
- $p < .005 \rightarrow$  „significant“
- $.005 < p < .05 \rightarrow$  „suggestive“
- $p > .05 \rightarrow$  „not significant“

Congratulations:  
You hacked your results to  
significance!

Now: What's the probability that a  
significant  $p$ -value indicates a true  
effect?





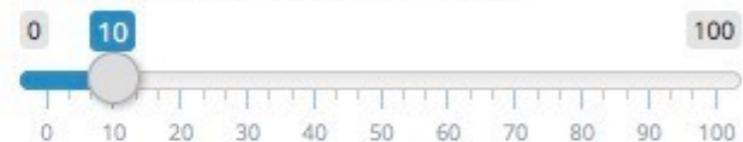


# When does a significant $p$ -value indicate a true effect?

## Understanding the Positive Predictive Value (PPV) of a $p$ -value

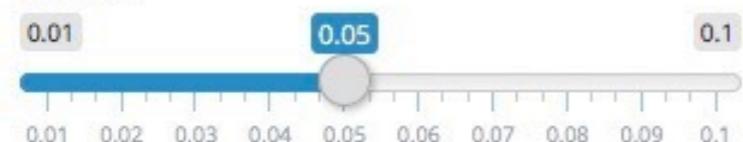
Across all investigated hypotheses: What % of them is actually true?

% of a priori true hypotheses:



What is your Type I error ( $\alpha$ ; typically 5%)?

$\alpha$  level



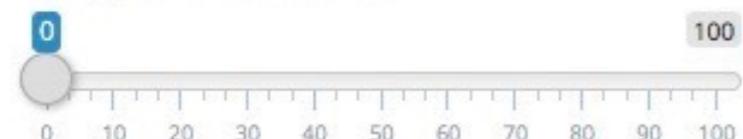
On what power level are the studies conducted?

Power



% of studies that report a significant result, although it's not ?

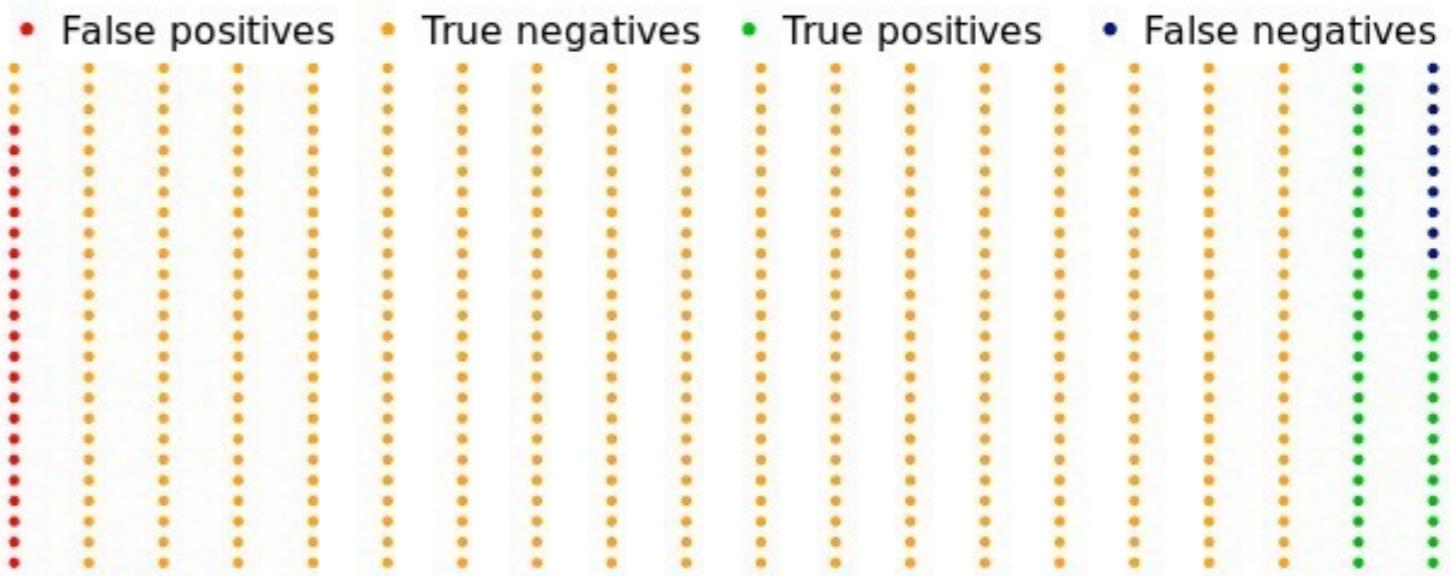
% of p-hacked studies



**Positive predictive rate (PPV):** 64% of claimed findings are true

**False discovery rate (FDR):** 36% of claimed findings are false

If we consider all findings, it looks like this (each point is one study):



If we consider **only the significant** findings, the ratio of true to false positives looks like this:



Practice with the PPV app! <http://shinyapps.org/apps/PPV/>

\* My lawyer told me to show that.

# Disclaimer\*

- $p$ -hacking increases the false positive rate
- $p$ -hacking „renders the reported  $p$ -values essentially uninterpretable“ (ASA statement)
- $p$ -hacking is ethically wrong and violates rules of good scientific practice
- If you  $p$ -hack systematically:
  - many of your research results will simply be wrong (depending on the prior probability of your hypotheses)
  - consequentially, your research won't replicate
- Every time you  $p$ -hack, you waste public money, you waste participants' time, you bias the literature, and **a kitten dies\*\*.**



\*\* If your research is about feline drug development

Gaming the system

or

Producing real knowledge?



„This was the largest audience ever to witness an inauguration, period.“

Press Secretary Sean Spicer



**Larry Parnell**  
@larry\_parnell

Starting immediately, all public-facing documents from USDA ARS will not be released. Disclaimer: On this forum, I don't represent my Agency



Folgen

post-truth! Fake-News!

# Trump's budget director pick: "Do we really need government-funded research at all?"

Mick Mulvaney suggested Zika science is uncertain, so we shouldn't bother to fund it.

Updated by Julia Belluz | @juliaoftoronto | julia.belluz@voxmedia.com | Jan 24, 2017, 11:07am EST

 TWEET

 SHARE



<http://www.vox.com/science-and-health/2016/12/budget-director-research-science-mulvaney>

Culture

## Donald Trump set to 'eliminate arts funding programs', cutting off NPR and PBS

He is notoriously not fond of broadcast media

Christopher Hooton | @christophhooton | 2 months ago |  319

 Gefällt

Click to follow  
The Independent Culture



<http://www.independent.co.uk/arts-entertainment/donald-trump-budget-cuts-arts-humanities-nea-neh-npr-cbs-president-a7536741.html>

# Norms

Communality

Open sharing

Universalism

Evaluate research on own merit

Disinterestedness

Motivated by knowledge and discovery

Organized skepticism

Consider all new evidence, even against  
one's prior work

Quality

# Counternorms

Secrecy

Closed

Particularism

Evaluate research by reputation

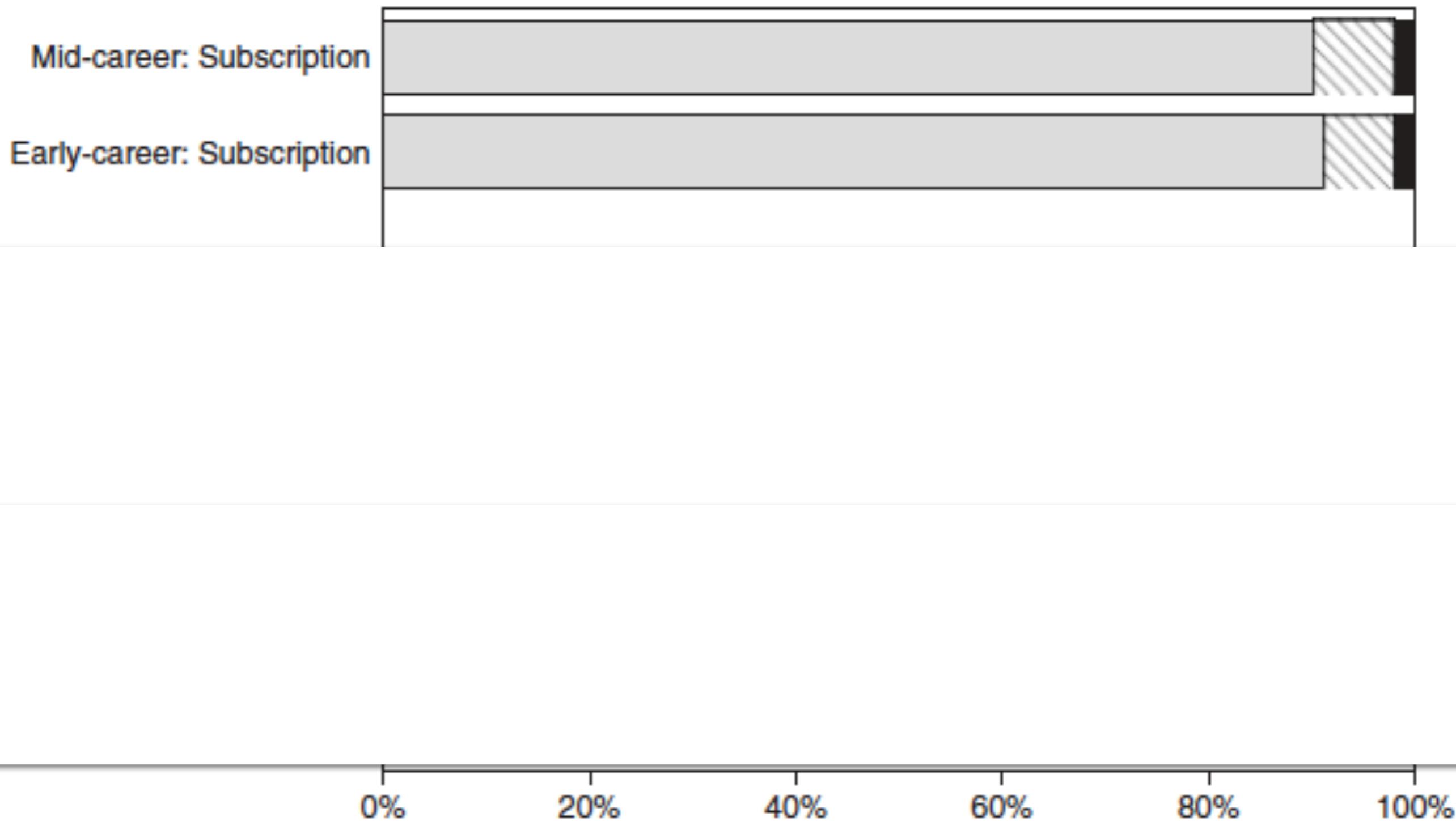
Self-interestedness

Treat science as a competition

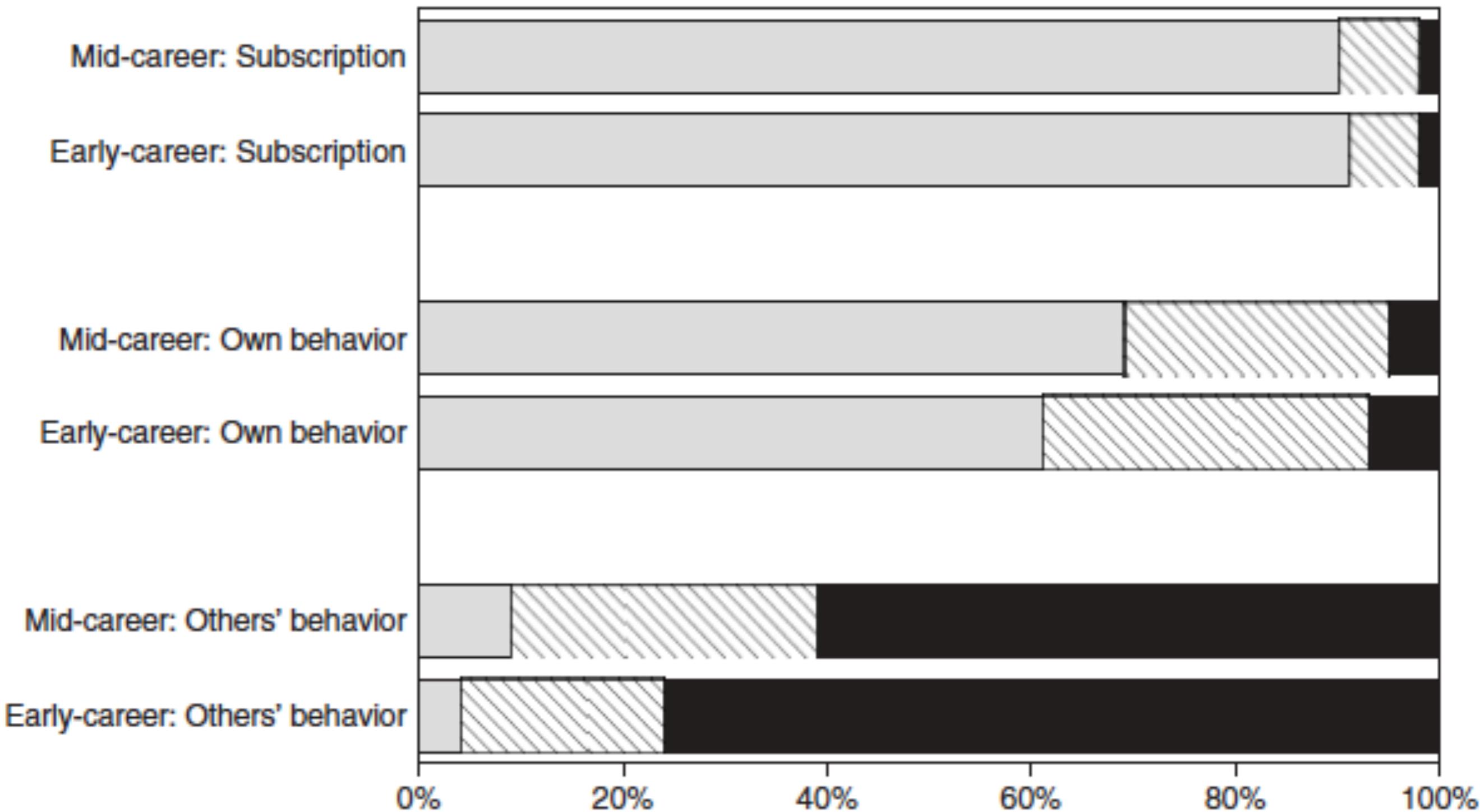
Organized dogmatism

Invest career promoting one's own  
theories, findings

Quantity



**FIG. 3. Norm versus Counternorm Scores: Percent with Norm > Counternorm (dotted), Norm = Counternorm (striped), Norm < Counternorm (solid).**



**FIG. 3. Norm versus Counternorm Scores: Percent with Norm > Counternorm (dotted), Norm = Counternorm (striped), Norm < Counternorm (solid).**