

# Replicability I

Fred LaPolla, MLS

Alisa Surkis, PhD, MLS

NYU Health Sciences Library



# Topics

- What is reproducibility
- The role of reproducibility today
- Study Design



# Outcomes: Students will be able to...

- Identify key types of reproducibility/replicability
- Explain why our understanding of concepts like "power" and "alpha" impact replicability
- Discuss how blinding and randomization can impact outcomes



# What does "reproducibility" mean to you?

# Terms

- Reproducibility
- Replicability
- Repeatability



# Science Translational Medicine

What does research reproducibility mean?

Steven N. Goodman\*, Daniele Fanelli and John P. A. Ioannidis

arXiv.org > cs > arXiv:1802.03311

Computer Science > Digital Libraries

## Terminologies for Reproducible Research

Lorena A. Barba

nature|methods

Correspondence | Published: 30 July 2015

Clarifying the terminology that describes scientific reproducibility

Ron S Kenett & Galit Shmueli



frontiers  
in Neuroinformatics

## Reproducibility vs. Replicability: A Brief History of a Confused Terminology

Hans E. Plesser<sup>1,2\*</sup>



Cold  
Spring  
Harbor  
Laboratory

bioRxiv

## A statistical definition for reproducibility and replicability

Prasad Patil, Roger D. Peng, Jeffrey Leek

NYU Langone  
Health





THE WALL STREET JOURNAL

World U.S. Politics Economy Business Tech Markets Opinion **Life & Arts** Real Estate

IDEAS | THE SATURDAY ESSAY

## The Breakdown in Biomedical Research

Contaminated samples, faulty studies and inadequate training have created a laboratories and industry, slowing the quest for new treatments and cures

npr WNYC RADIO news arts & life music programs shop

## HIDDEN BRAIN

A CONVERSATION ABOUT LIFE'S UNSEEN PATTERNS

28:17

### When Great Minds Think Unlike: Inside Science's 'Replication Crisis'

May 24, 2016 - 12:10 AM ET

## The New York Times

### *Many Psychology Findings Not as Strong as Claimed, Study Says*

## Sunday Review

### Why Do So Many Studies Fail to Replicate?

The Economist World politics Business & finance Economics Science & technology Culture

Unreliable research

### Trouble at the lab

Scientists like to think of science as self-correcting. To an alarming degree, it is not

Oct 19th 2013 | From the print edition

Like 11k Tweet 1,227



Lisson Ford



# National Academies of Sciences 2019 Report:

- Reproducibility: Obtaining consistent results from the same data, computation and analysis (**computational reproducibility**)
- Replicability: Obtaining consistent results across studies on the same question
- Generalizability: The extent that one study applies to other contexts





## Ask:

- If someone uses my code to process and analyze my raw data, will I get an identical answer?
- If I repeat the same experiment, will I get a result that is consistent with my original result?
- If someone else tries to replicate my experiment, will they get a result that is consistent with my result?
- Will someone else replicating the experiment draw a conclusion that is consistent with the original?



# For this class: Reproducibility vs Replicability

## Reproducibility

- Typically expect bitwise agreement
- Exact reproduction of result does not guarantee it's correct

## Replicability

- Do not expect every study to replicate
- Lack of replication does not necessarily indicate flaws in the experimental process

## Factors Limiting Replicability: The Good

- Complexity of the system under study
- Intrinsic variation in nature
- Variables outside the scope of current scientific knowledge
- Limitations of current technologies
- Prior probability of the scientific hypothesis (i.e. unexpected results more likely not to replicate)

# Factors Limiting Replicability: The Bad

- Poor study design
- Poor execution
- Misuse/understanding of statistics
- Researcher bias
- Publication bias

# Factors Limiting Replicability: The Ugly

- Fraud
- Relatively rare



NATURE | COMMENT

## A long journey to reproducible results

Gordon J. Lithgow, Monica Driscoll & Patrick Phillips

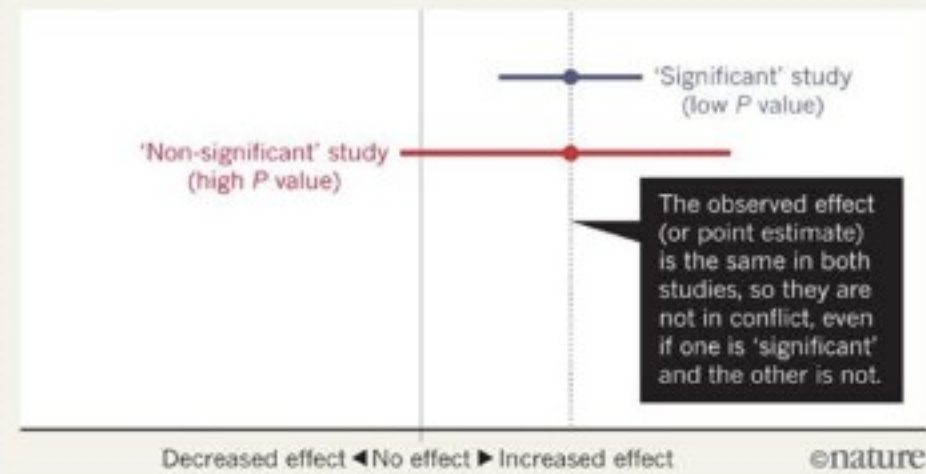
# A 2017 Replicability Study

- Aimed to learn about aging in worms
- Lab technique variability slowed progress
- Found that there were actually cohorts of worms that partition into short and long life



### BEWARE FALSE CONCLUSIONS

Studies currently dubbed 'statistically significant' and 'statistically non-significant' need not be contradictory, and such designations might cause genuine effects to be dismissed.



Source: V. Amrhein et al.

# No litmus test for replicability

- What is being replicated: direction of effect? magnitude?
- Significance can be problematic
- Inherent uncertainty in statistics



NATURE | NEWS



## Over half of psychology studies fail reproducibility test

Largest replication study to date casts doubt on many published positive results.

### Response:



## A Bayesian Perspective on the Reproducibility Project: Psychology

Alexander Etz , Joachim Vandekerckhove

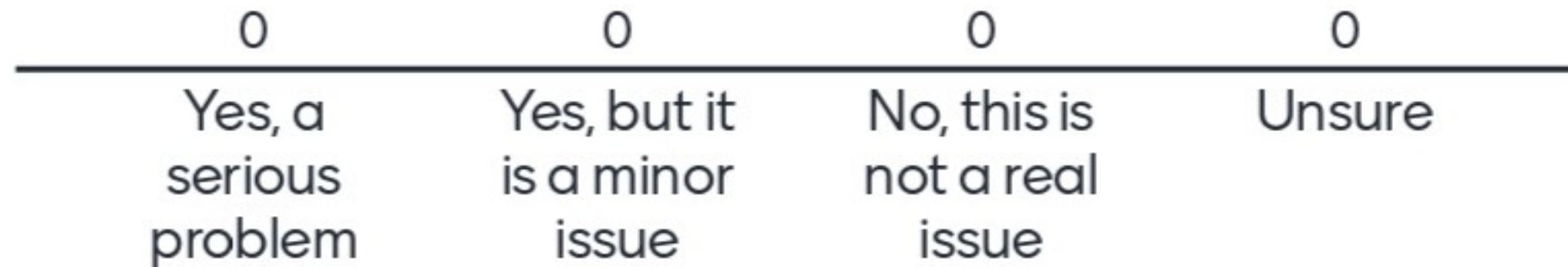
Published: February 26, 2016 • DOI: 10.1371/journal.pone.0149794

“...analysis revealed no obvious inconsistencies between the original and replication results.”

“...apparent failure of the Reproducibility Project to replicate many target effects can be adequately explained by overestimation of effect sizes... due to **small sample sizes** and **publication bias** in the psychological literature.



# Do you feel that there is a reproducibility crisis?





# Why reproducibility matters







**National Institutes  
of Health**

## NIH Guidelines as of 2016

- Premise
- Rigorous study design
- Relevant biological variables
- Authentication of resources



# The NIH Names Factors Contributing to Poor Reproducibility

- Poor training in experimental design
- Focus on headline grabbing statements
- Lack of detail in publications
- Publication bias
- Misinterpretation of hypothesis-generating research



# Review

## Reproducibility in Science Improving the Standard for Basic and Preclinical Research

C. Glenn Begley, John P.A. Ioannidis

## 2015 Review Article

**Table 1. Examples of Some Reported Reproducibility Concerns in Preclinical Research**

Author	Field
Ioannidis et al (2009) <sup>12</sup>	Microarray data
Baggerly et al (2009) <sup>13</sup>	Microarray data
Sena et al (2010) <sup>14</sup>	Stroke animal studies
Prinz (2011) <sup>1</sup>	General biology
Begley & Ellis (2012) <sup>2</sup>	Oncology
Nekrutenko & Taylor (2012) <sup>15</sup>	NGS data access
Perrin (2014) <sup>16</sup>	Mouse, in-vivo
Tsilidis et al (2013) <sup>17</sup>	Neurological studies
Lazic & Essioux (2013) <sup>18</sup>	Mouse VPA model
Haibe-Kains et al (2013) <sup>19</sup>	Genomics/cell line analysis
Witwer (2013) <sup>20</sup>	Microarray data
Elliott et al (2006) <sup>21</sup>	Commercial antibodies
Prassas et al (2013) <sup>22</sup>	Commercial ELISA
Stodden et al (2013) <sup>23</sup>	Journals
Baker et al (2014) <sup>24</sup>	Journals
Vaux (2012) <sup>25</sup>	Journals

ALS indicates amyotrophic lateral sclerosis; MIAME, minimum information about a microarray experiment; NGS, next generation sequencing.

**Table 2. Additional Basic Science Fields Where Concerns Regarding Reproducibility Have Been Raised**

Discipline	Issues Raised
Neuroscience	Low statistical power; small sample size
Pharmacology	Lack of training, lack of statistical power, blinding, hypothesis, requisite PK studies, randomization, dose-response, controls, prospective plan, validation, independent replication, and selection of doses that are not tolerable in humans
Genomics/bioinformatics	Irreproducibility of high-profile studies
Stem cell biology	Lack of reliable, quality data
Oncology, in vitro testing	Use of clinically unachievable concentrations
Chemistry lead-discovery	Artifacts; false positives and negatives
Computational biology	10 common errors
Pathology/Biomarkers	Biospecimen quality
Organizational psychology	Suppression of negative studies
Observational research	Q/52 hypotheses confirmed in randomized trials

<https://www.ahajournals.org/doi/pdf/10.1161/CIRCRESAHA.114.303819>





# nature

International weekly journal of science

## Drug development: Raise standards for preclinical cancer research

C. Glenn Begley & Lee M. Ellis

Scientists in hematology and oncology departments at Amgen tried to confirm findings from 53 “landmark” studies



# How many of the 53 landmark studies were confirmed?





# Preclinical Research

- In 6 (11%) studies, findings were confirmed
- 25% were consistent enough to continue research (per the standards of Amgen)
- <https://www.nature.com/articles/483531a>



# What do scientists think?

NATURE | NEWS FEATURE

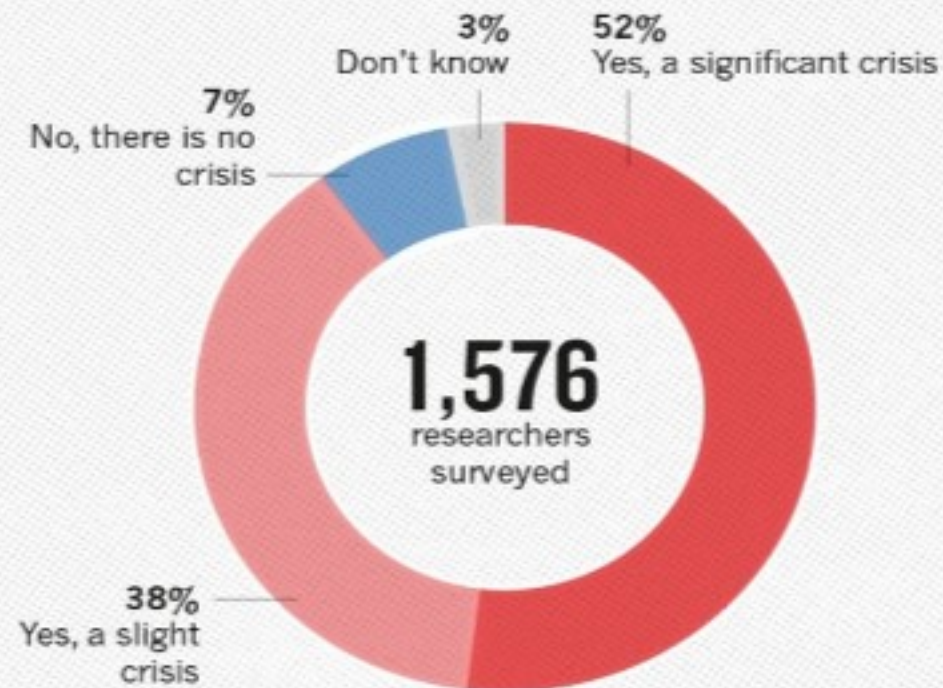
## 1,500 scientists lift the lid on reproducibility

Survey sheds light on the 'crisis' rocking research.

Monya Baker

25 May 2016 |

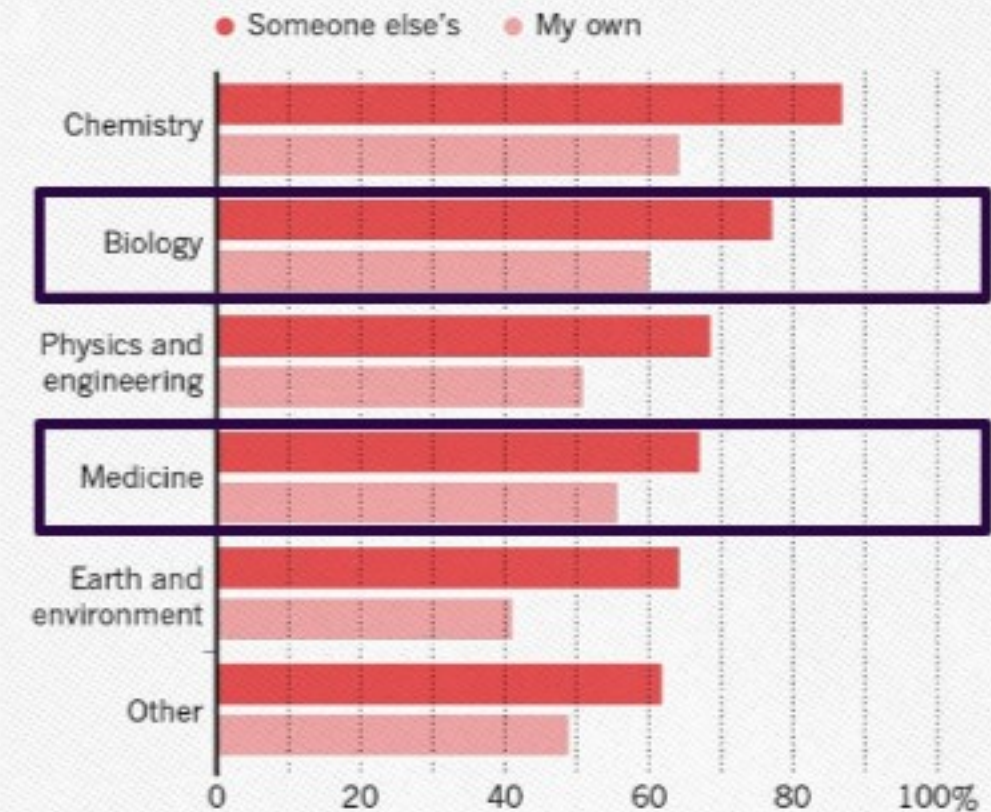
### IS THERE A REPRODUCIBILITY CRISIS?



©nature

### HAVE YOU FAILED TO REPRODUCE AN EXPERIMENT?

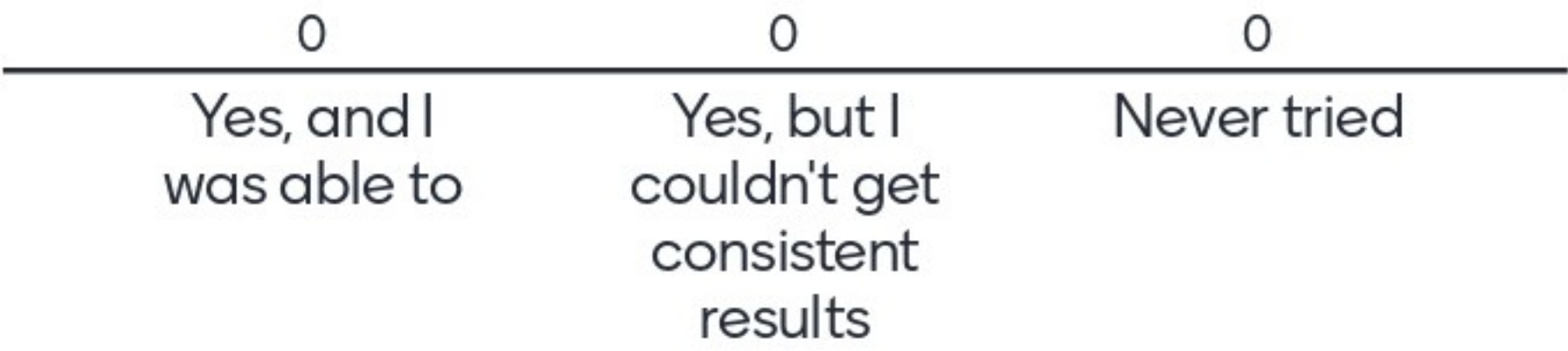
Most scientists have experienced failure to reproduce results.



<http://www.nature.com/news/1-500-scientists-lift-the-lid-on-reproducibility-1.19970>



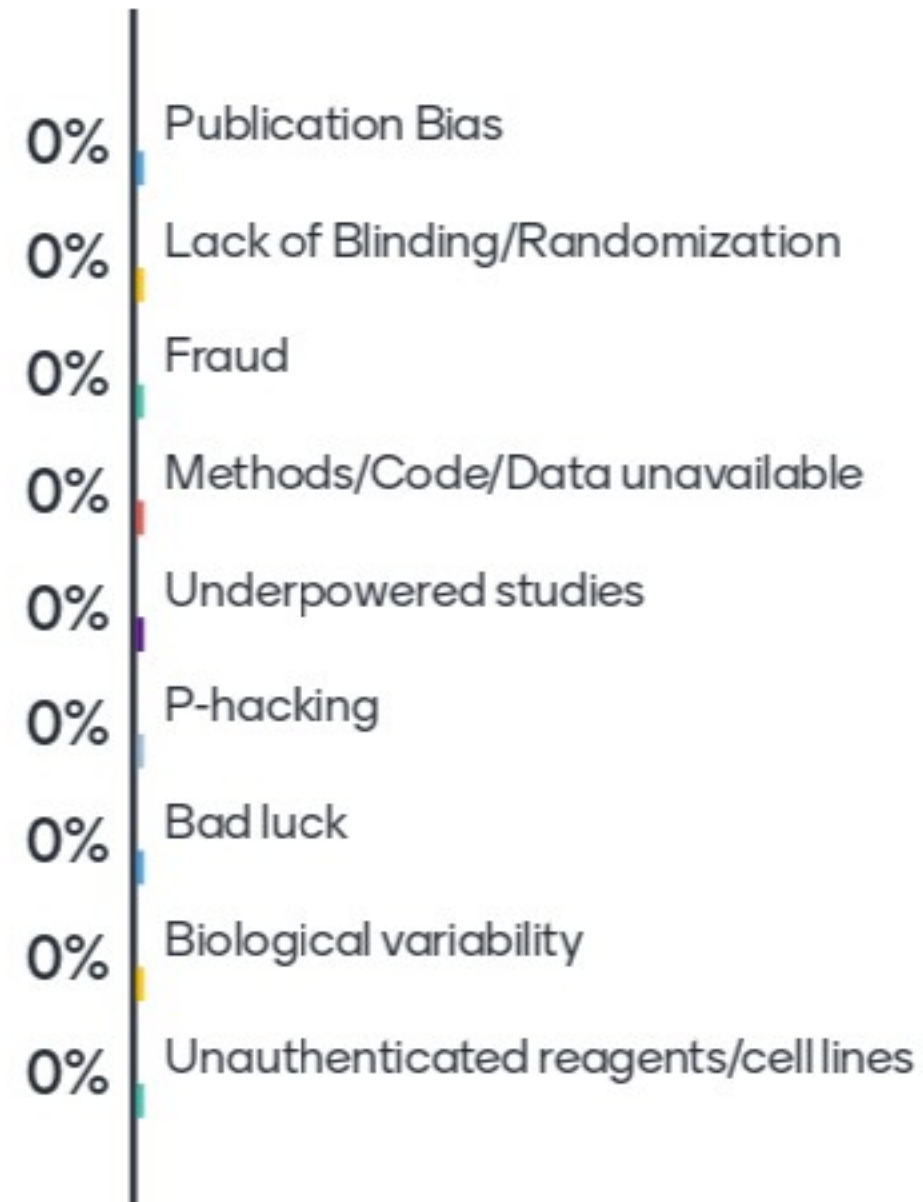
# Have you ever tried to replicate someone else's study?





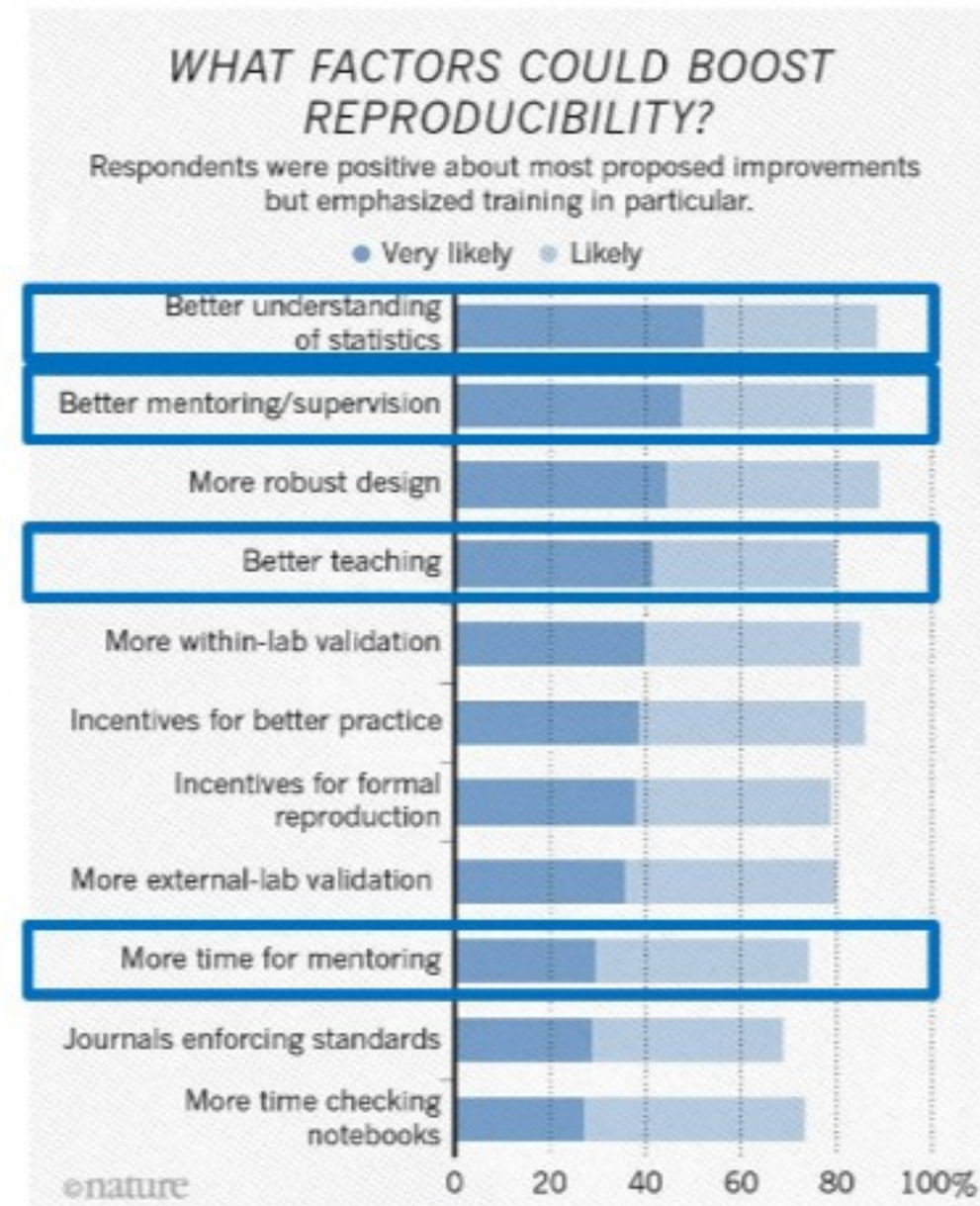
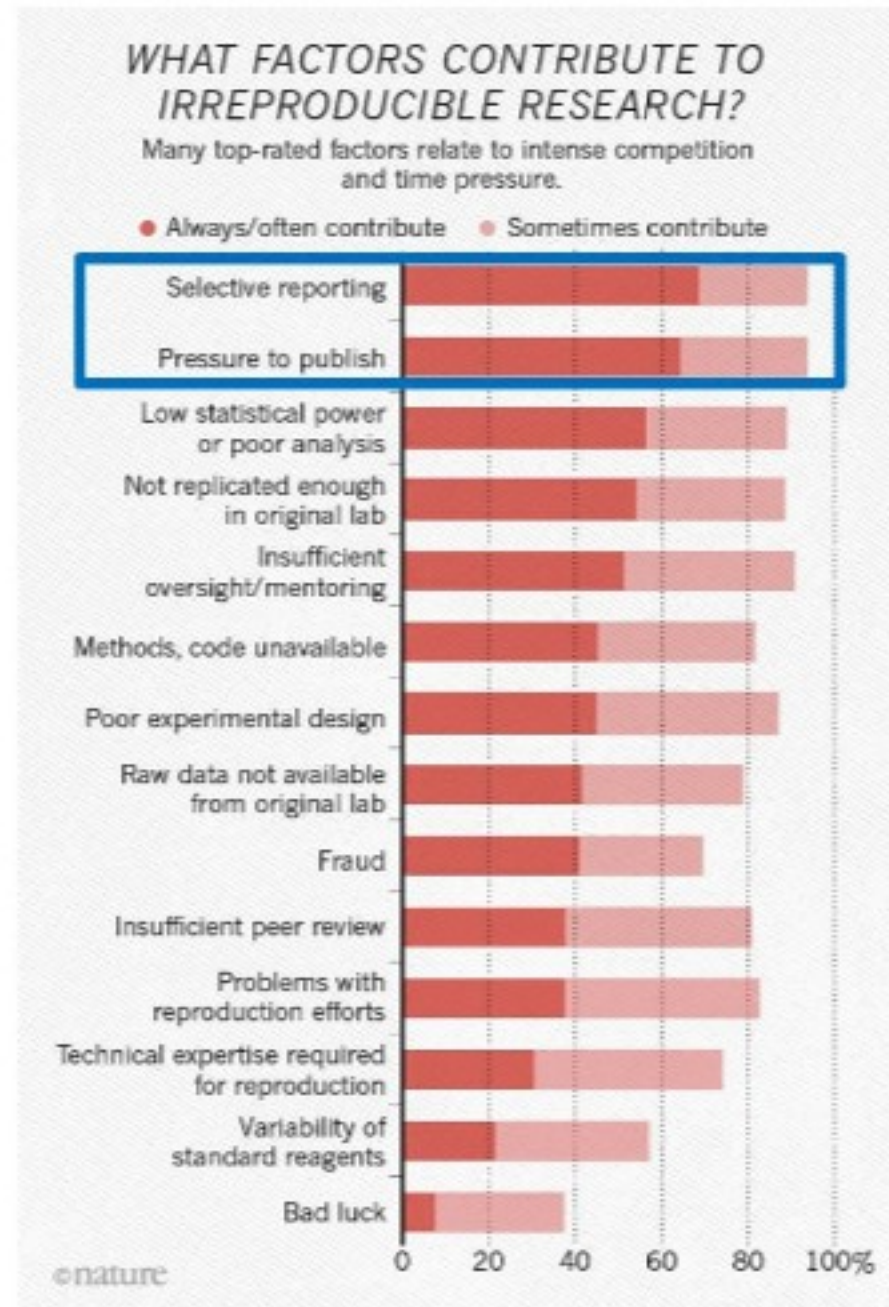
If you tried and were unable to replicate: what were the barriers and what was learned?

Assign points to each based on how big a factor you think it is in the current reproducibility/replicability crisis





# What do scientists think?



<http://www.nature.com/news/1-500-scientists-lift-the-lid-on-reproducibility-1.19970>



## The Economics of Reproducibility in Preclinical Research PLoS Biology, 2015

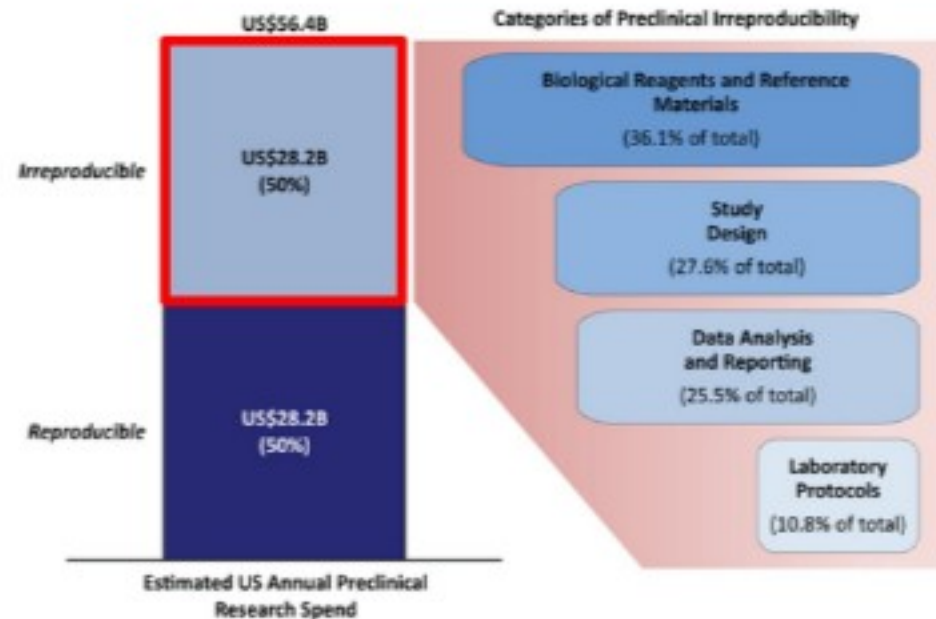


Fig 2. Estimated US preclinical research spend and categories of errors that contribute to irreproducibility.

# Estimated \$28 Billion on Irreproducible Research

- Looked at study design, biological resources, protocols and analysis
- Estimated an upper and lower bound of impact on reproducibility
- Estimate based on \$56 Billion total on preclinical research

<https://journals.plos.org/plosbiology/article?id=10.1371/journal.pbio.1002165>





# Study Design

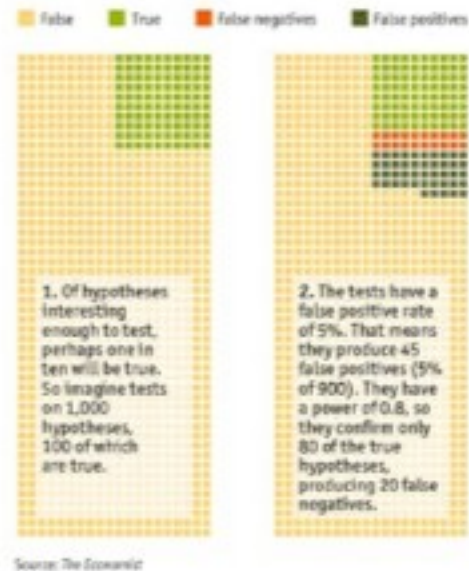


## Power!

- Statistical power is the probability a test will discover differences when they exist
- Higher power lowers the probability of Type II Error (Power is the probability of avoiding Type II error)
- 0.8 is a commonly used benchmark
- Power is influenced by sample size, effect size and variance
- Larger samples are needed to detect smaller effects







# A thought experiment


- 100/1000 hypotheses are true
- If Power is 0.8, we expect to find 80 true positives will be detected
- If  $\alpha = .05$ , we can expect to find  $0.5 \times 900 = 45$  false positives
- Of the 125 positive findings, 36% (45) may be false

nature  
REVIEWS

NEUROSCIENCE

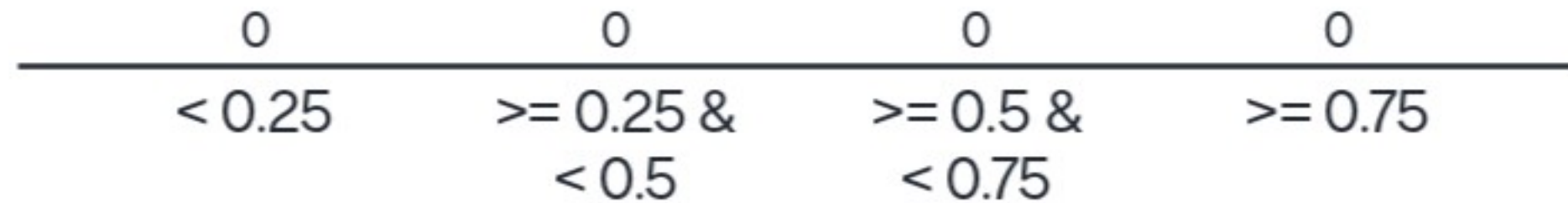
# Power failure: why small sample size undermines the reliability of neuroscience

Katherine S. Button, John P. A. Ioannidis, Claire Mokrysz, Brian A. Nosek, Jonathan Flint, Emma S. J. Robinson & Marcus R. Munafò 

*Nature Reviews Neuroscience* **14**, 365–376 (2013) | [Download Citation](#) 

A 2013 study looked at the power levels in neuroscience studies

# In what range do you think the median statistical power fell?

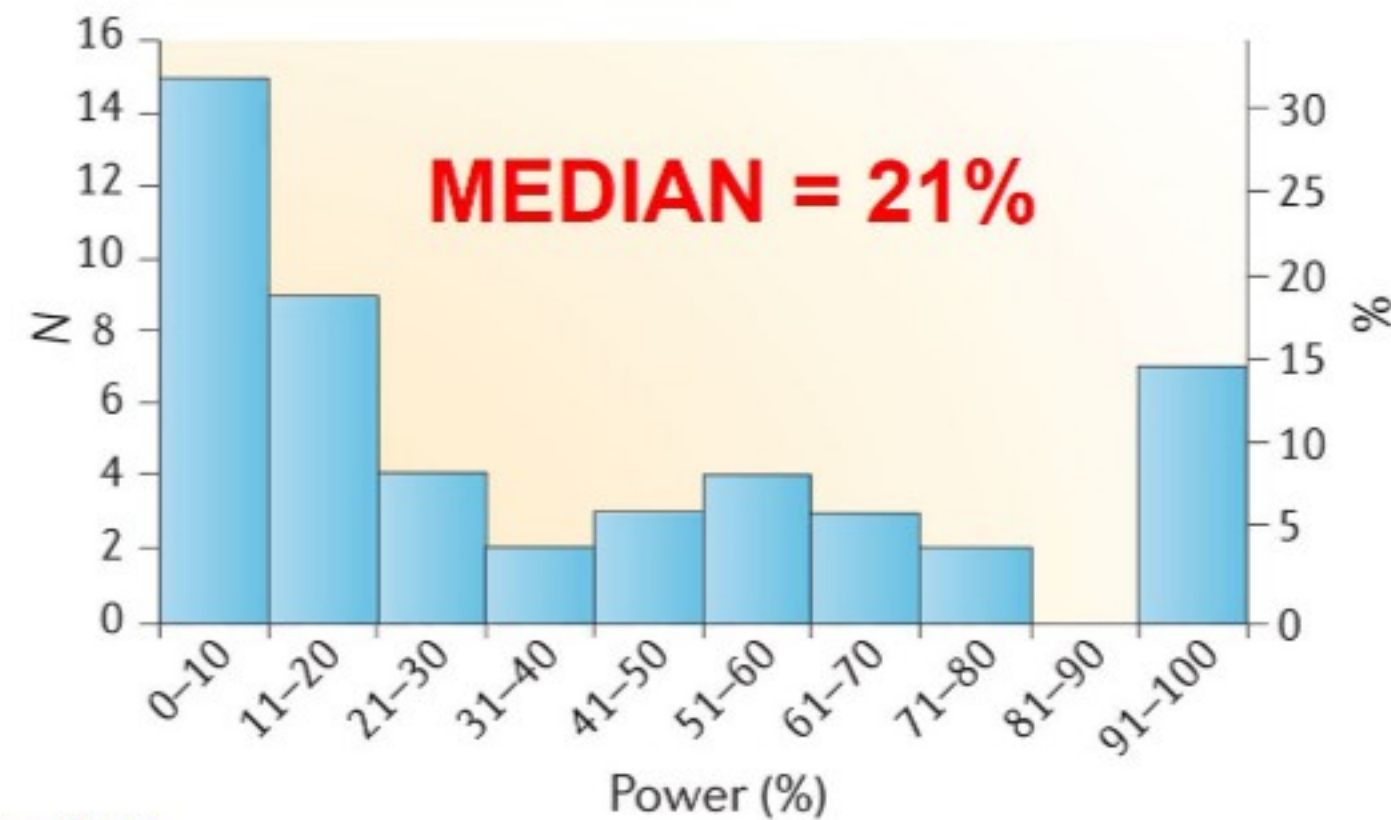




# Power failure: why small sample size undermines the reliability of neuroscience

Katherine S. Button, John P. A. Ioannidis, Claire Mokrysz, Brian A. Nosek, Jonathan Flint, Emma S. J. Robinson & Marcus R. Munafò

*Nature Reviews Neuroscience* **14**, 365–376 (2013) | [Download Citation](#)



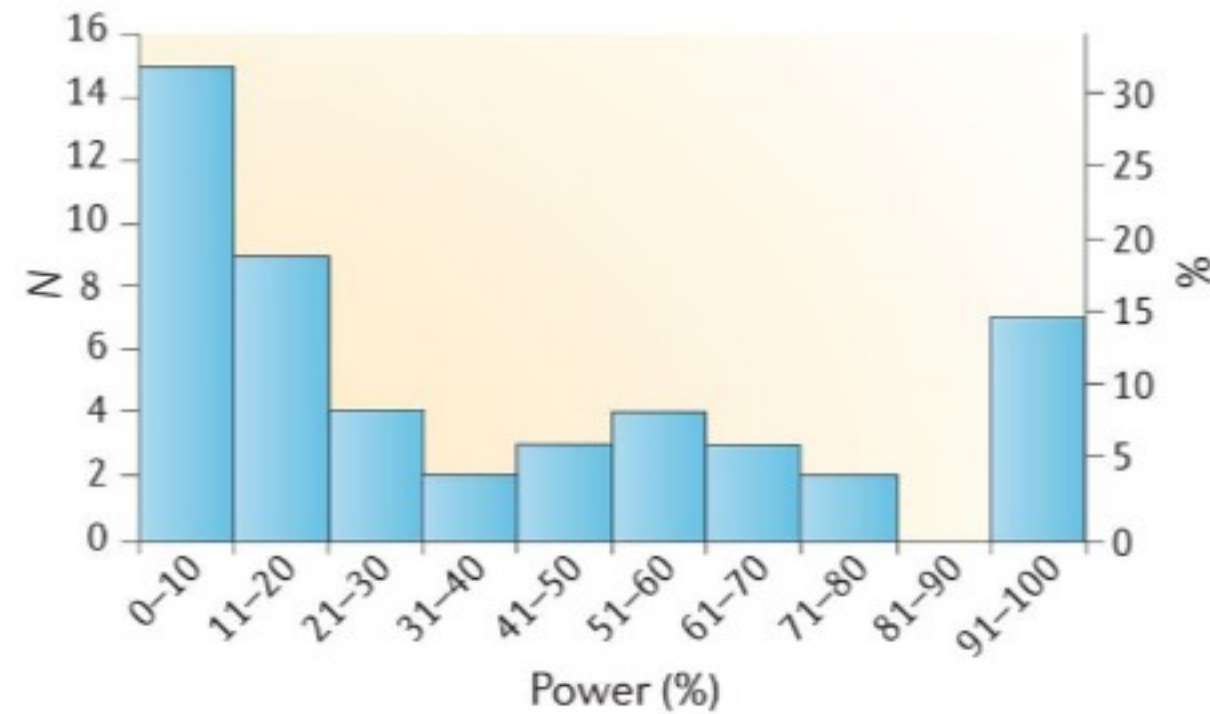
<https://www.nature.com/articles/nrn3475>

A power of .21 greatly increases the likelihood of type II error





For power of 0.20, how many times out of 100 would we expect to detect a difference between groups when one is present?



- 100/1000 hypotheses are true
- Power = 0.2

- $0.2 * 100 =$  power to detect 20 True Positives
- $900 * .05 =$  possibility to miss 45 False Positives
- $45/65 = 69\%$  of findings may be False Positives



# "Winner's Curse"

- For small study, a smaller actual effect size will not reach statistical significance
- Therefore, low powered study only large effects will be "significant"
- "Lucky" scientist who discovers effect in a small study is cursed by finding an inflated effect



# Ways to determine a meaningful effect size

- Look at effect in pilots
- Clinical/subject expertise
- Variation
- Consult with a statistician (NYU Biostatistics Resource)



# Bias



# In the context of science, what is bias?

# Bias:

The **systematic** introduction of error



# COGNITIVE FALLACIES IN RESEARCH



## HYPOTHESIS MYOPIA

Collecting evidence to support a hypothesis, not looking for evidence against it, and ignoring other explanations.



## TEXAS SHARPSHOOTER

Seizing on random patterns in the data and mistaking them for interesting findings.



## ASYMMETRIC ATTENTION

Rigorously checking unexpected results, but giving expected ones a free pass.



## JUST-SO STORYTELLING

Finding stories after the fact to rationalize whatever the results turn out to be.





# DEBIASING TECHNIQUES



## DEVIL'S ADVOCACY

Explicitly consider alternative hypotheses — then test them out head-to-head.



## PRE-COMMITMENT

Publicly declare a data collection and analysis plan before starting the study.



## TEAM OF RIVALS

Invite your academic adversaries to collaborate with you on a study.



## BLIND DATA ANALYSIS

Analyse data that look real but are not exactly what you collected — and then lift the blind.



*Lack of blinding of outcome assessors in animal model experiments implies risk of observer bias*

## Blinding:

- 2014 meta-analysis looked at impact of not blinding outcome assessors on estimates of intervention effects in 10 animal studies (2,450 animals)
- Found that unblinded studies exaggerated odds ratios of effect by 59%
- Important to note this does not imply bad intent
- <https://www.sciencedirect.com/science/article/abs/pii/S0895435614001577>



# Randomization

- 2014 systematic review found failure to randomize leads to overestimation of treatment effect
- <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0098856>





Do you feel failure to blind or randomize may result in bias? Why or why not?



## Blinding and Randomization Not Happening

- 2019 study of 574 papers found 56% reported if randomization happened and blinding of outcomes assessors happened in 31%
- <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0215221>
- 2014 study of 2280 papers found 25% randomized, 15-24% blinded
- <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0089981>



# In studies you have done, have the experiments

0	0	0	0
Been both blinded and randomized	Blinded to outcomes only	Randomized subjects only	Neither blinded nor randomized



**Review!**



# List factors that limit replicability



# Why does low power limit replicability?

# Why might a lack of blinding and randomization limit replicability?

# Homework

Watch this video: <https://www.youtube.com/watch?v=a4fUU85ABwc>

|

Write a paragraph (~200 words) addressing the following:

Have labs you've worked in generally employed randomization? Do you think randomization are generally employed in your area of study? If you have employed randomization in a study, did you consider all the factors discussed in the video?

|

Due in Brightspace by 9:00am 2/12/2020



# Bibliography

1. National Academies of Sciences, Engineering, and Medicine. 2019. *Reproducibility and Replicability in Science*. Washington, DC: The National Academies Press.
2. Lithgow GJ, Driscoll M, Phillips P. A long journey to reproducible results. *Nature*. 2017;548(7668):387-8.
3. Amrhein V, Greenland S, McShane B. Scientists rise up against statistical significance. *Nature*. 2019 Mar;567(7748):305-307.
4. Baker, M. Over half of psychology studies fail reproducibility test. *Nature*. 2015 August 27.
5. Etz A, Vandekerckhove J. A Bayesian Perspective on the Reproducibility Project: Psychology. *PLOS ONE*. 2016;11(2):e0149794.
6. Stark PB. Before reproducibility must come preproducibility. *Nature*. 2018;557(7707):613.
7. Freedman LP, Cockburn IM, Simcoe TS. The Economics of Reproducibility in Preclinical Research. *PLOS Biology*. 2015;13(6):e1002165.
8. Begley CG, Ioannidis JP. Reproducibility in science: improving the standard for basic and preclinical research. *Circulation research*. 2015;116(1):116-26.
9. Begley CG, Ellis LM. Drug development: Raise standards for preclinical cancer research. *Nature*. 2012 Mar 28;483(7391):531-3. doi: 10.1038/483531a.
10. Prinz F, Schlange T, Asadullah K. Believe it or not: how much can we rely on published data on potential drug targets? *Nature Reviews Drug Discovery*. 2011;10:712.
11. Baker M. 1,500 scientists lift the lid on reproducibility. *Nature*. 2016;533(7604):452-4.
12. Nuzzo R. How scientists fool themselves - and how they can stop. *Nature*. 2015 Oct 8;526(7572):182-5.
13. Avey MT, Moher D, Sullivan KJ, Fergusson D, Griffin G, et al. (2016) The Devil Is in the Details: Incomplete Reporting in Preclinical Animal Research. *PLOS ONE* 11(11): e0166733.
14. Ioannidis JPA. Why Most Published Research Findings Are False. *PLOS Medicine*. 2005;2(8):e124.
15. Button KS, Ioannidis JPA, Mokrysz C, Nosek BA, Flint J, Robinson ESJ, et al. Power failure: why small sample size undermines the reliability of neuroscience. *Nature Reviews Neuroscience*. 2013;14:365.
16. Bello S, Krogsbøll LT, Gruber J, Zhao ZJ, Fischer D, Hróbjartsson A. Lack of blinding of outcome assessors in animal model experiments implies risk of observer bias. *J Clin Epidemiol*. 2014 Sep;67(9):973-83.
17. Bebert V, Luyten D, Heard K. Emergency medicine animal research: does use of randomization and blinding affect the results? *Acad Emerg Med*. 2003 Jun;10(6):684-7.
18. Hirst JA, Howick J, Aronson JK, Roberts N, Perera R, et al. (2014) The Need for Randomization in Animal Trials: An Overview of Systematic Reviews. *PLOS ONE* 9(6): e98856.
19. Kilkenny C, Parsons N, Kadyszewski E, Festing MFW, Cuthill IC, Fry D, et al. Survey of the Quality of Experimental Design, Statistical Analysis and Reporting of Research Using Animals. *PLOS ONE*. 2009;4(11):e7824.
20. van Luijk J, Bakker B, Rovers MM, Ritskes-Hoitinga M, de Vries RB, Leenaars M. Systematic reviews of animal studies; missing link in translational research? *PLoS One*. 2014 Mar 26;9(3):e89981.
21. Fergusson DA, Avey MT, Barron CC, Bocock M, Bieffer KE, et al. (2019) Reporting preclinical anesthesia study (REPEAT): Evaluating the quality of reporting in the preclinical anesthesiology literature. *PLOS ONE* 14(5): e0215221.

