## Robust tests of theory with randomly sampled experiments

Joachim Vandekerckhove, Beth Baribault, Chris Donkin, Daniel R. Little, Jennifer Trueblood, Zita Oravecz, Don van Ravenzwaaij, Corey White

University of California, Irvine

and various other institutions of great renown

## Reproducibility

2

- ... is widely considered to be in crisis (Baker, 2016)

- ... is widely considered to be in crisis (Baker, 2016)
- ... is the topic of a National Academy of Sciences review (Cicerone, 2015)

- ... is widely considered to be in crisis (Baker, 2016)
- ... is the topic of a National Academy of Sciences review (Cicerone, 2015)
- ... is now one of the fundamental review criteria for all grants submitted to the National Institutes of Health

- ... is widely considered to be in crisis (Baker, 2016)
- ... is the topic of a National Academy of Sciences review (Cicerone, 2015)
- ... is now one of the fundamental review criteria for all grants submitted to the National Institutes of Health
- ... has been discussed in Congressional hearings

- ... is widely considered to be in crisis (Baker, 2016)
- ... is the topic of a National Academy of Sciences review (Cicerone, 2015)
- ... is now one of the fundamental review criteria for all grants submitted to the National Institutes of Health
- ... has been discussed in Congressional hearings
- ... "threatens the entire biomedical research enterprise" (Rosenblatt, 2016).

There is broad support for reproducibility-boosting corrective measures (Baker, 2016):

There is broad support for reproducibility-boosting corrective measures (Baker, 2016):

- Better statistics

There is broad support for reproducibility-boosting corrective measures (Baker, 2016):

- Better statistics
- Internal and external validation of research results

There is broad support for reproducibility-boosting corrective measures (Baker, 2016):

- Better statistics
- Internal and external validation of research results
- Robust experimental designs

## Robustness

"Well, but ...."

"Well, but ..."

- "... your study took place in Amsterdam; ours in Nijmegen."

"Well, but ..."

- "... your study took place in Amsterdam; ours in Nijmegen."
- "... your study took place in 2017; ours in 1999."

"Well, but ..."

- "... your study took place in Amsterdam; ours in Nijmegen."

- "... your study took place in 2017; ours in 1999."

I would question the societal value of scientific claims that do not generalize beyond 1999 Nijmegen.

The concept of hidden moderators is as old as science.

The concept of hidden moderators is as old as science.

nial cinnabar: for though in our furnaces it hath been very fuccefsfully made, yet not only we have afterwards failed of making it, but we have feen much more expert chymifts, and who, becaufe of the high value they do (not undefervedly) place upon that medicine, imploy themfelves oftener than we in making it, divers times unfuccefsfully attempt the preparing it. And it may be perhaps also from fome diverfity either in antimonies or irons, that eminent chymifts have (as we have obferved) oftenfailed in their endeavours to make the flarry regulus of *Mars* and antimony. Infomuch that divers artifts fondly believe and teach (what our experience will not permit us to allow) that there is a certain refpect to times and conficultations requifite to the producing of this (1 confefs admirable) body. Upon which fubject I muft

Boyle (1772), on the suspected importance of astrological time to chemical purification of antimony.

## Terminology

 Replicability: How can we ensure that other people can repeat the study procedures we used?

- Replicability: How can we ensure that other people can repeat the study procedures we used?
- Reproducibility: How can we ensure that other people can repeat the analysis we performed?

- Replicability: How can we ensure that other people can repeat the study procedures we used?
- Reproducibility: How can we ensure that other people can repeat the analysis we performed?
- Robustness: How can we ensure that other people obtain the same findings?

- Replicability: How can we ensure that other people can repeat the study procedures we used?
- Reproducibility: How can we ensure that other people can repeat the analysis we performed?
- Robustness: How can we ensure that other people obtain the same findings?
- Generalizability: How can we ensure that our results translate to very different contexts?

- Replicability: How can we ensure that other people can repeat the study procedures we used? → Open materials
- Reproducibility: How can we ensure that other people can repeat the analysis we performed? → Code hygiene; open data
- Robustness: How can we ensure that other people obtain the same findings?  $\rightarrow$  ? Robust design?
- Generalizability: How can we ensure that our results translate to very different contexts? → ? Higher abstraction?

 Robust statistics [have] good performance for data drawn from a wide range of probability distributions, especially for distributions that are not normal

- Robust statistics [have] good performance for data drawn from a wide range of probability distributions, especially for distributions that are not normal
- Robust programming is a style of programming that focuses on handling ... unexpected actions. It requires code to handle these ... actions gracefully

- Robust statistics [have] good performance for data drawn from a wide range of probability distributions, especially for distributions that are not normal
- Robust programming is a style of programming that focuses on handling ... unexpected actions. It requires code to handle these ... actions gracefully
- [In engineering,] Robustness is the state where the technology, product, or process performance is minimally sensitive to factors causing variability (either in the manufacturing or users environment)

 A low bar for robustness would seem to be "extend to a future repetition of an identical procedure"

- A low bar for robustness would seem to be "extend to a future repetition of an identical procedure"
- A slightly higher (but still low) bar would be "extend to a procedure that only differs in irrelevant ways" (e.g., lab)

- A low bar for robustness would seem to be "extend to a future repetition of an identical procedure"
- A slightly higher (but still low) bar would be "extend to a procedure that only differs in irrelevant ways" (e.g., lab)
- One way of thinking about this is to imagine a population of largely interchangeable experiments

- A low bar for robustness would seem to be "extend to a future repetition of an identical procedure"
- A slightly higher (but still low) bar would be "extend to a procedure that only differs in irrelevant ways" (e.g., lab)
- One way of thinking about this is to imagine a population of largely interchangeable experiments
- We have methods to license statements about populations: sampling
# **Randomized design**

Three independent groups invoke similar sounding metaphors:

- Universe of generalization (Cronbach et al., 1963)
- Constraints of generality (Simons et al., 2017)
- Boundary of meaning (Kenett & Rubinstein, 2017)

Three independent groups invoke similar sounding metaphors:

- Universe of generalization (Cronbach et al., 1963)
- Constraints of generality (Simons et al., 2017)
- Boundary of meaning (Kenett & Rubinstein, 2017)

These metaphors imply the existence of some spatially-arranged population of possible experiments where the effect holds.

Three independent groups invoke similar sounding metaphors:

- Universe of generalization (Cronbach et al., 1963)
- Constraints of generality (Simons et al., 2017)
- Boundary of meaning (Kenett & Rubinstein, 2017)

These metaphors imply the existence of some spatially-arranged population of possible experiments where the effect holds.

To generalize across this population of experiments, we could sample from it.

 Reuss et al. (2015) report on a subliminal cueing effect



Time course of one trial.

- Reuss et al. (2015) report on a subliminal cueing effect
- The cue b or v tells the participant to favor speed or accuracy in the task



Time course of one trial.

- Reuss et al. (2015) report on a subliminal cueing effect
- The cue b or v tells the participant to favor speed or accuracy in the task
- When the cue is masked, it becomes subliminal



Time course of one trial.

- Reuss et al. (2015) report on a subliminal cueing effect
- The cue b or v tells the participant to favor speed or accuracy in the task
- When the cue is masked, it becomes subliminal
- Even so, participants changed their speed-accuracy trade-off



Time course of one trial.

- Characters used

- Characters used
- Color of the stimuli

- Characters used
- Color of the stimuli
- Brightness of the stimuli

- Characters used
- Color of the stimuli
- Brightness of the stimuli
- Identity of the participants

- Characters used
- Color of the stimuli
- Brightness of the stimuli
- Identity of the participants

- ...

These "theoretically inert features" of an experiment are largely ones that could go lost in the translation from a verbal claim to an experimental design.

These "theoretically inert features" of an experiment are largely ones that could go lost in the translation from a verbal claim to an experimental design.

These independent-variable dimensions span a space that we'll call the method space.

These "theoretically inert features" of an experiment are largely ones that could go lost in the translation from a verbal claim to an experimental design.

These independent-variable dimensions span a space that we'll call the method space.

Theoretical statements often correspond to functions over the method space (e.g., an effect size is nonzero over some region).

These "theoretically inert features" of an experiment are largely ones that could go lost in the translation from a verbal claim to an experimental design.

These independent-variable dimensions span a space that we'll call the method space.

Theoretical statements often correspond to functions over the method space (e.g., an effect size is nonzero over some region).

Experiments are points or small areas in that space.

Now we sample from these points and conduct one micro-experiment (n = 32) for each location



Each micro-experiment has 32 trials in a  $2 \times 2$  design followed by 32 trials of a cue identification task.

- Now we sample from these points and conduct one micro-experiment (n = 32) for each location
- Often these are exchangeable for our theoretical purposes



Each micro-experiment has 32 trials in a  $2 \times 2$  design followed by 32 trials of a cue identification task.

- Now we sample from these points and conduct one micro-experiment (n = 32) for each location
- Often these are exchangeable for our theoretical purposes
- Then jointly analyze in a "planned" meta-analysis



Each micro-experiment has 32 trials in a  $2 \times 2$  design followed by 32 trials of a cue identification task.

- Now we sample from these points and conduct one micro-experiment (n = 32) for each location
- Often these are exchangeable for our theoretical purposes
- Then jointly analyze in a "planned" meta-analysis
- We should expect to see some level-2 variability



Each micro-experiment has 32 trials in a  $2 \times 2$  design followed by 32 trials of a cue identification task.

 Focusing only on unmasked trials, the effect seems reasonably robust



- Focusing only on unmasked trials, the effect seems reasonably robust
- Participants are 10-15mm more accurate after an accuracy prompt



- Focusing only on unmasked trials, the effect seems reasonably robust
- Participants are 10-15mm more accurate after an accuracy prompt
- 78% of micro-experiments show evidence for the effect



- Focusing only on unmasked trials, the effect seems reasonably robust
- Participants are 10-15mm more accurate after an accuracy prompt
- 78% of micro-experiments show evidence for the effect
- This is not surprising since this is just the positive control



 Focusing only on masked trials, the effect seems fickle



## **Planned meta-analysis**

- Focusing only on masked trials, the effect seems fickle
- 75% of micro-experiments are more consistent with an effect of 0mm than 10mm



Histogram of effect sizes  $\Delta$ mm (masked condition).

## **Planned meta-analysis**

- Focusing only on masked trials, the effect seems fickle
- 75% of micro-experiments are more consistent with an effect of 0mm than 10mm
- Lots of level-2 variability



- Focusing only on masked trials, the effect seems fickle
- 75% of micro-experiments are more consistent with an effect of 0mm than 10mm
- Lots of level-2 variability
- Does not support a robust subliminal cueing effect



Histogram of effect sizes  $\Delta$ mm (masked condition).

- Focusing only on masked trials, the effect seems fickle
- 75% of micro-experiments are more consistent with an effect of 0mm than 10mm
- Lots of level-2 variability
- Does not support a robust subliminal cueing effect
- But there might be an effect under limited conditions



Histogram of effect sizes  $\Delta$ mm (masked condition).

 As a manipulation check, we measured cue identifiability under the unique conditions of each micro-experiment

- As a manipulation check, we measured cue identifiability under the unique conditions of each micro-experiment
- We can now introduce cue detection performance as a covariate in a meta-regression

- As a manipulation check, we measured cue identifiability under the unique conditions of each micro-experiment
- We can now introduce cue detection performance as a covariate in a meta-regression
- The effect size function is flat in the subliminal range



Effect sizes ( $\Delta$ mm) as a function of detection performance.
As robustness check, include target color as a covariate

- As robustness check, include target color as a covariate
- Flat function indicates robustness to target color



Effect sizes  $(\Delta mm)$  as a function of target color, split by detection performance.

- As robustness check, include target color as a covariate
- Flat function indicates robustness to target color
- Conditional on supraliminality, the subliminal effect seems robustly present



Effect sizes ( $\Delta$ mm) as a function of target color, split by detection performance.

- As robustness check, include target color as a covariate
- Flat function indicates robustness to target color
- Conditional on supraliminality, the subliminal effect seems robustly present
- Conditional on subliminality, the subliminal effect seems robustly absent



Effect sizes ( $\Delta$ mm) as a function of target color, split by detection performance.

 permits statements at a higher level of abstraction: the existence of an effect in a well-defined and formalized universe of intended generalization (e.g., the varying design choices that other researchers might have made in this study)

- permits statements at a higher level of abstraction: the existence of an effect in a well-defined and formalized universe of intended generalization (e.g., the varying design choices that other researchers might have made in this study)
- allows for defensive design: a design strategy that targets replicability, reproducibility, robustness, and generalizability

- permits statements at a higher level of abstraction: the existence of an effect in a well-defined and formalized universe of intended generalization (e.g., the varying design choices that other researchers might have made in this study)
- allows for defensive design: a design strategy that targets replicability, reproducibility, robustness, and generalizability
- turns out to be feasible and relatively cheap

- permits statements at a higher level of abstraction: the existence of an effect in a well-defined and formalized universe of intended generalization (e.g., the varying design choices that other researchers might have made in this study)
- allows for defensive design: a design strategy that targets replicability, reproducibility, robustness, and generalizability
- turns out to be feasible and relatively cheap
- doesn't require very complicated statistical methods (yet)

# Discussion

 I think most real effects are probably robust



What I think real effects behave like.

- I think most real effects are probably robust
- The Many Labs projects (e.g., Klein et al., 2014) lend little support to the idea that effects routinely wink in and out of detectability



What I think real effects behave like.

- I think most real effects are probably robust
- The Many Labs projects (e.g., Klein et al., 2014) lend little support to the idea that effects routinely wink in and out of detectability
- Only robust claims have predictive power, which presumably aids practical applicability



What I think real effects behave like.

 I suspect many social scientists believe effects are often fickle



What I think many social scientists believe.

 I suspect many social scientists believe effects are often fickle

... making it important—and a bit of an art—to conduct very finely tuned experiments



What I think many social scientists believe.

- I suspect many social scientists believe effects are often fickle
  - ... making it important—and a bit of an art—to conduct very finely tuned experiments
- Valid or not, this mindset invites irreproducibility



What I think many social scientists believe.

- I suspect many social scientists believe effects are often fickle
  - ... making it important—and a bit of an art—to conduct very finely tuned experiments
- Valid or not, this mindset invites irreproducibility
- Fickle effects are systematically hard to distinguish from flukes



What I think many social scientists believe.



These two conceptions of real effects probably call for different kinds of empirical programs – strategies with different balances of exploration vs. exploitation.

More generally I think there exist methods that are optimized for sensitivity, robustness, and generalizability, that could be used more in psychology and social sciences.

More generally I think there exist methods that are optimized for sensitivity, robustness, and generalizability, that could be used more in psychology and social sciences.

Techniques like adding variability to designs to improve robustness and generalizability come directly out of the statistical / psychometrical literature on random effects. More generally I think there exist methods that are optimized for sensitivity, robustness, and generalizability, that could be used more in psychology and social sciences.

Techniques like adding variability to designs to improve robustness and generalizability come directly out of the statistical / psychometrical literature on random effects.

These techniques are not difficult to implement.

More generally I think there exist methods that are optimized for sensitivity, robustness, and generalizability, that could be used more in psychology and social sciences.

Techniques like adding variability to designs to improve robustness and generalizability come directly out of the statistical / psychometrical literature on random effects.

These techniques are not difficult to implement.

 $\sim$  fin  $\sim$ 

### Some (current and recent) benefactors



## JOHN TEMPLETON

#### FOUNDATION



This is joint work with Beth Baribault, Chris Donkin, Daniel Little, Don van Ravenzwaaij, Jennifer Trueblood, Corey White, Zita Oravecz, and Paul de Boeck.



Baker, M. (2016). Is there a reproducibility crisis? Nature, 533, 453–455.
Baribault, B., Donkin, C., Little, D. R., Trueblood, J. S., Oravecz, Z., van Ravenzwaaij, D., et al. (2018). Metastudies for robust tests of theory. Proceedings of the National Academy of Sciences, 201708285.

Boyle, R. I. (1772). The works of the Honourable Robert Boyle, In six volumes, To which is prefixed the life of the author (Vol. 2). J. and F. Rivington.
Cicerone, R. (2015). Research reproducibility, replicability, reliability. National Academy of Sciences (NAS) Annual Meeting.

Cronbach, L. J., Rajaratnam, N., & Gleser, G. C. (1963). Theory of generalizability: A liberalization of reliability theory. *British Journal of Statistical Psychology*, 16(2), 137–163.

#### References ii

Kenett, R. S., & Rubinstein, A. (2017). A generalization approach to reproducibility claims.

(https://papers.ssrn.com/sol3/papers.cfm?abstract\_id=3035070)

- Klein, R. A., Ratliff, K. A., Vianello, M., Adams Jr, R. B., Bahník, Š., Bernstein, M. J., et al. (2014). Investigating variation in replicability: A "many labs" replication project. *Social psychology*, 45(3), 142.
- Reuss, H., Kiesel, A., & Kunde, W. (2015). Adjustments of response speed and accuracy to unconscious cues. *Cognition*, 134, 57–62.
- Rosenblatt, M. (2016). An incentive-based approach for improving data reproducibility. *Science Translational Medicine*, 8(336), 336ed5–336ed5.
  Simons, D. J., Shoda, Y., & Lindsay, D. S. (2017). Constraints on generality

(COG): A proposed addition to all empirical papers. *Perspectives on Psychological Science*, *12*(6), 1123–1128.