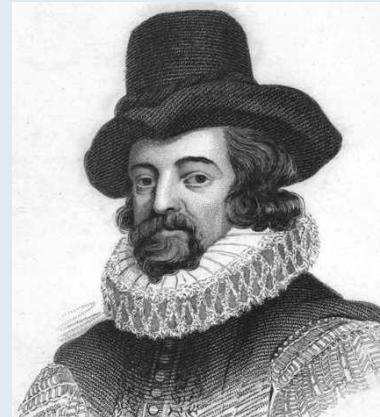


Psychological and sociological motivations and challenges for blind analyses

Rob MacCoun

*Scientia nihil aliud est quam
veritatis imago. (Science is but
an image of the truth.)*

— Francis Bacon



*Scientia nihil aliud est quam
veritatis imago. (Science is but
an image of the truth.)*

— Francis Bacon

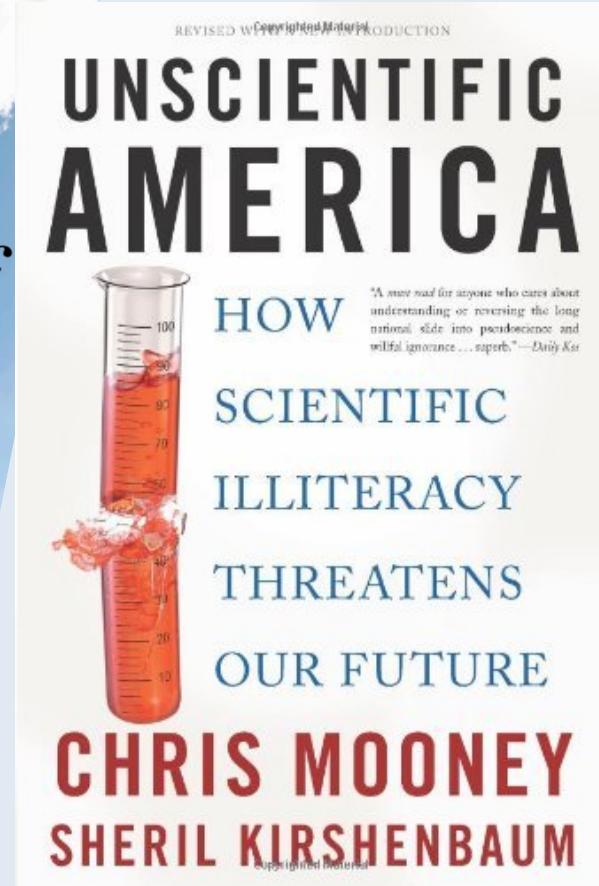
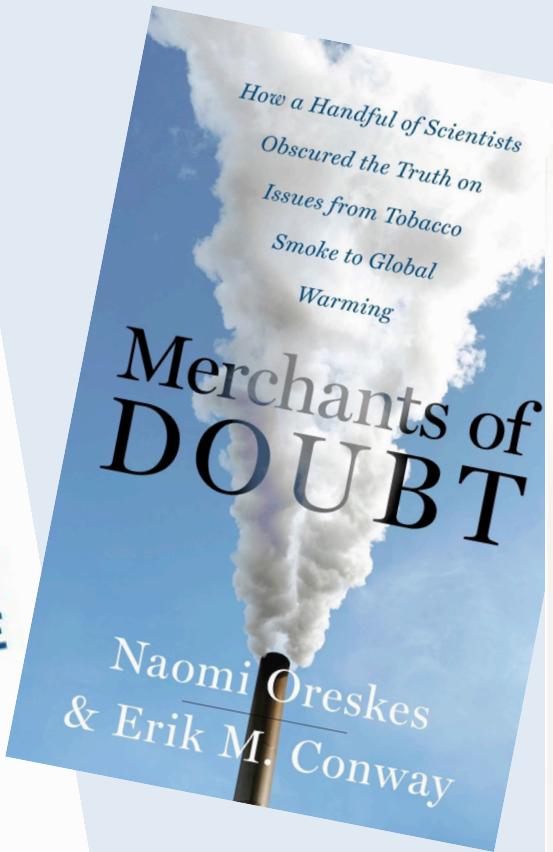
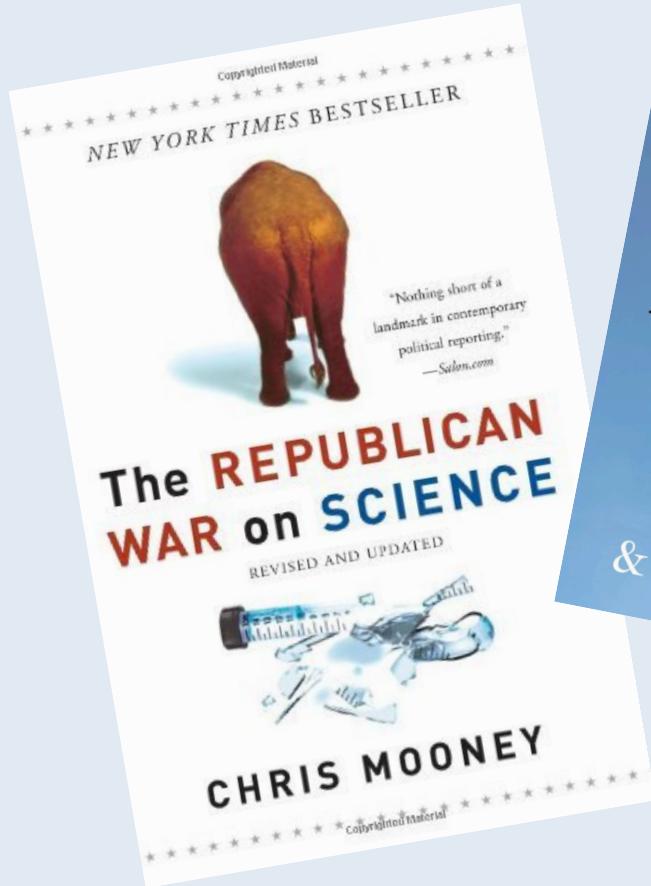


Back off man. I'm a scientist.

— Bill Murray, playing Dr. Peter Venkman, in the movie *Ghostbusters*.







SOCIAL SELECTION

*Popular articles
on social media*

A proposal to cure sloppy science

A debate this week on way quickly spread to social media. College London, called ‘Is science fix it?’, included claims practices, such as tweaking results seem significant, at com/4imbij). One suggestion to register their experiments with a journal before support. Audience publisher at the op-

NATURE.COM

For more on
popular papers:
go.nature.com/mjyia

SOCIAL SELECTION

*Popular topics
on social media*

Science reporting flaws exposed

Researchers often complain about inaccurate science stories in the popular press, but few air their grievances in a

The New York Times | <http://nyti.ms/1JlfMXV>

The Opinion Pages | OP-ED CONTRIBUTORS

What's Behind Big Science Frauds?

By ADAM MARCUS and IVAN ORANSKY MAY 22, 2015

NATURE.COM

For more on
popular papers:
go.nature.com/awqdox

Canada, posted on Facebook:

“Journalists and scientists should both strive for accuracy in their claims.”

Front Psychol. 6, 00988 (2015)

SCIENCE

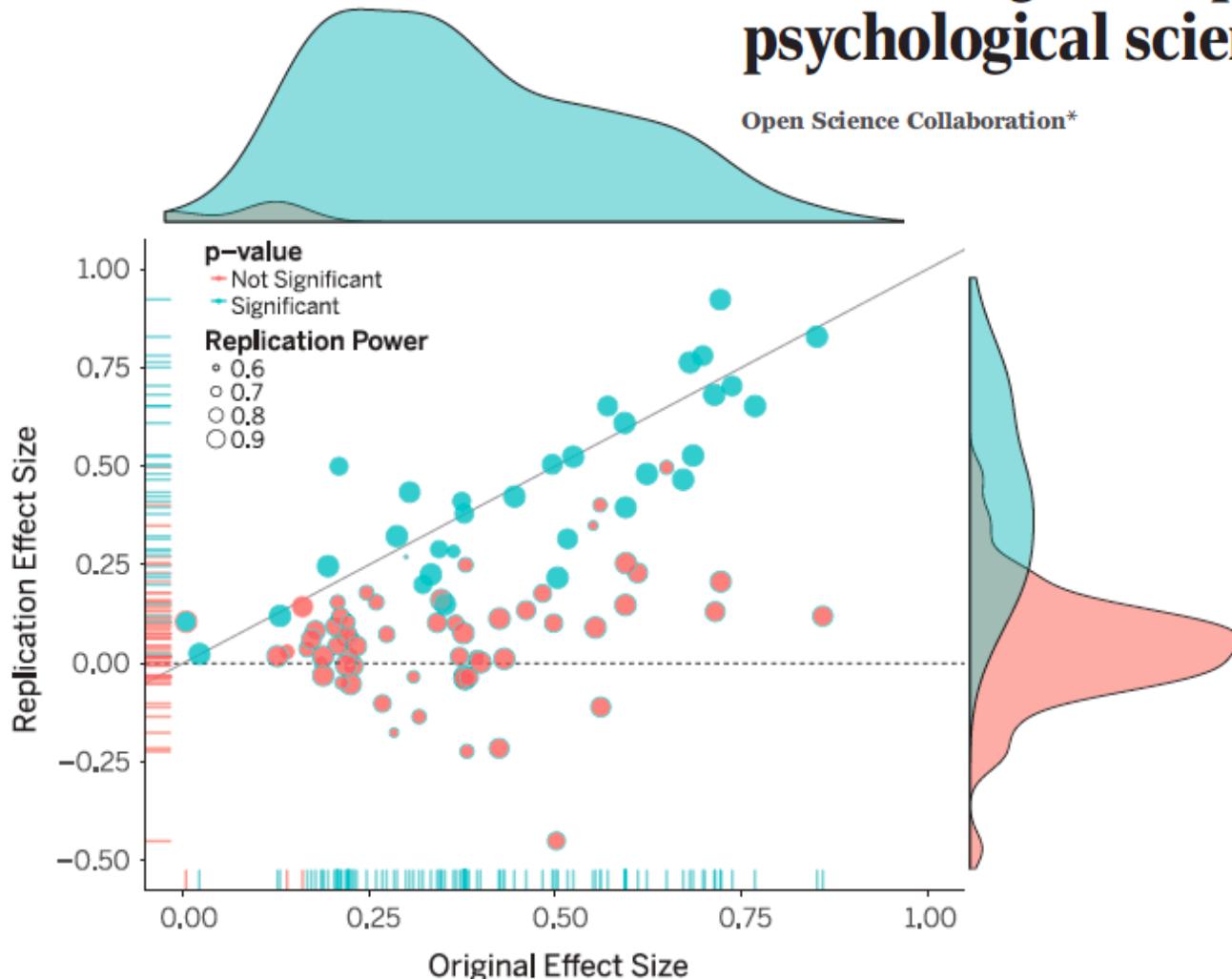
Many Psychology Findings Not as Strong as Claimed, Study Says

By BENEDICT CAREY AUG. 27, 2015

Estimating the reproducibility of psychological science

Open Science Collaboration*

Estimating the reproducibility of psychological science



Original study effect size versus replication effect size (correlation coefficients). Diagonal line represents replication effect size equal to original effect size. Dotted line represents replication effect size of 0. Points below the dotted line were effects in the opposite direction of the original. Density plots are separated by significant (blue) and nonsignificant (red) effects.

Rescuing US biomedical research from its systemic flaws

Bruce Alberts^a, Marc W. Kirschner^b, Shirley Tilghman^{c,1}, and Harold Varmus^d

^aDepartment of Biophysics and Biochemistry, University of California, San Francisco, CA 94158; ^bDepartment of Systems Biology, Harvard Medical School, Boston, MA 02115; ^cDepartment of Molecular Biology, Princeton University, Princeton, NJ 08540; and ^dNational Cancer Institute, Bethesda, MD 20892

NATURE | NEWS

Irreproducible biology research costs put at \$28 billion per year

Study calculates cost of flawed biomedical research in the United States.

Monya Baker

09 June 2015

Systematic Review of the Empirical Evidence of Study Publication Bias and Outcome Reporting Bias

Kerry Dwan^{1*}, Douglas G. Altman², Juan A. Arnaiz³, Jill Bloom⁴, An-Wen Chan⁵, Eugenia Cronin⁶, Evelyne Decullier⁷, Philippa J. Easterbrook⁸, Erik Von Elm^{9,10}, Carrol Gamble¹, Davina Ghersi¹¹, John P. A. Ioannidis^{12,13}, John Simes¹⁴, Paula R. Williamson¹

Annu. Rev. Psychol. 1998. 49:259–87

Copyright © 1998 by Annual Reviews Inc. All rights reserved

BIASES IN THE INTERPRETATION AND USE OF RESEARCH RESULTS

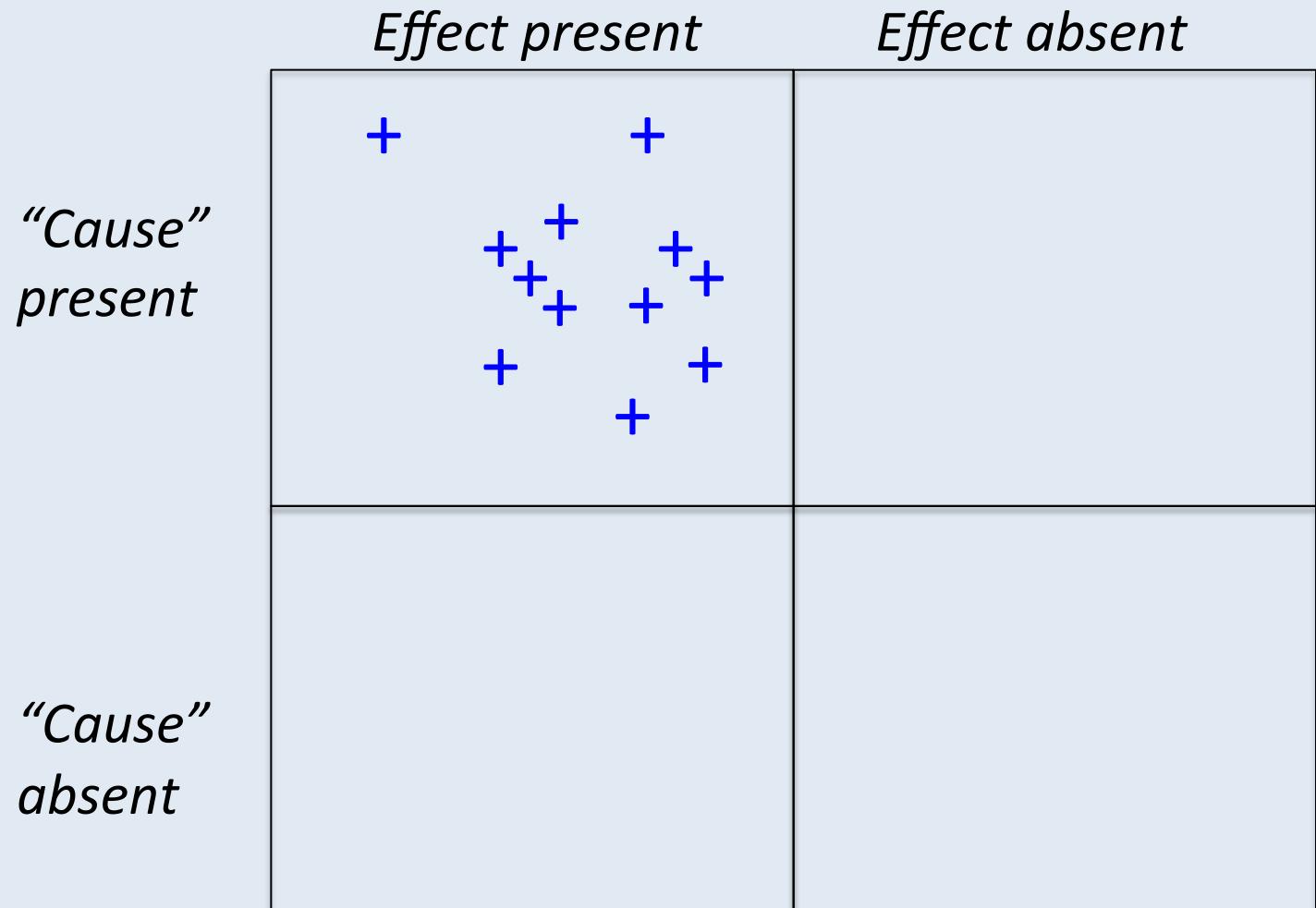
Robert J. MacCoun

Bias Prototypes

	<i>Intentional?</i>	<i>Motivated?</i>	<i>Justifiable?</i>
<i>Fraud</i>	Yes	Yes	No
<i>Advocacy</i>	Yes	Yes	Maybe
<i>Cold Bias</i>	No	No	No
<i>Hot Bias</i>	No	Yes	No
<i>Skepticism</i>	Maybe	Yes	Yes

MacCoun (1998, *Ann.Rev.Psy.*)

Confirmatory bias in evidence search



Confirmation bias in evidence evaluation: Lord et al (1979)

- Participants were either *pro or anti capital punishment*
- Read 1 of 4 versions of a research article:

	Found deterrent effect	Found no deterrent effect
Longitudinal		
Cross-sectional		

Lord et al (1979): Findings

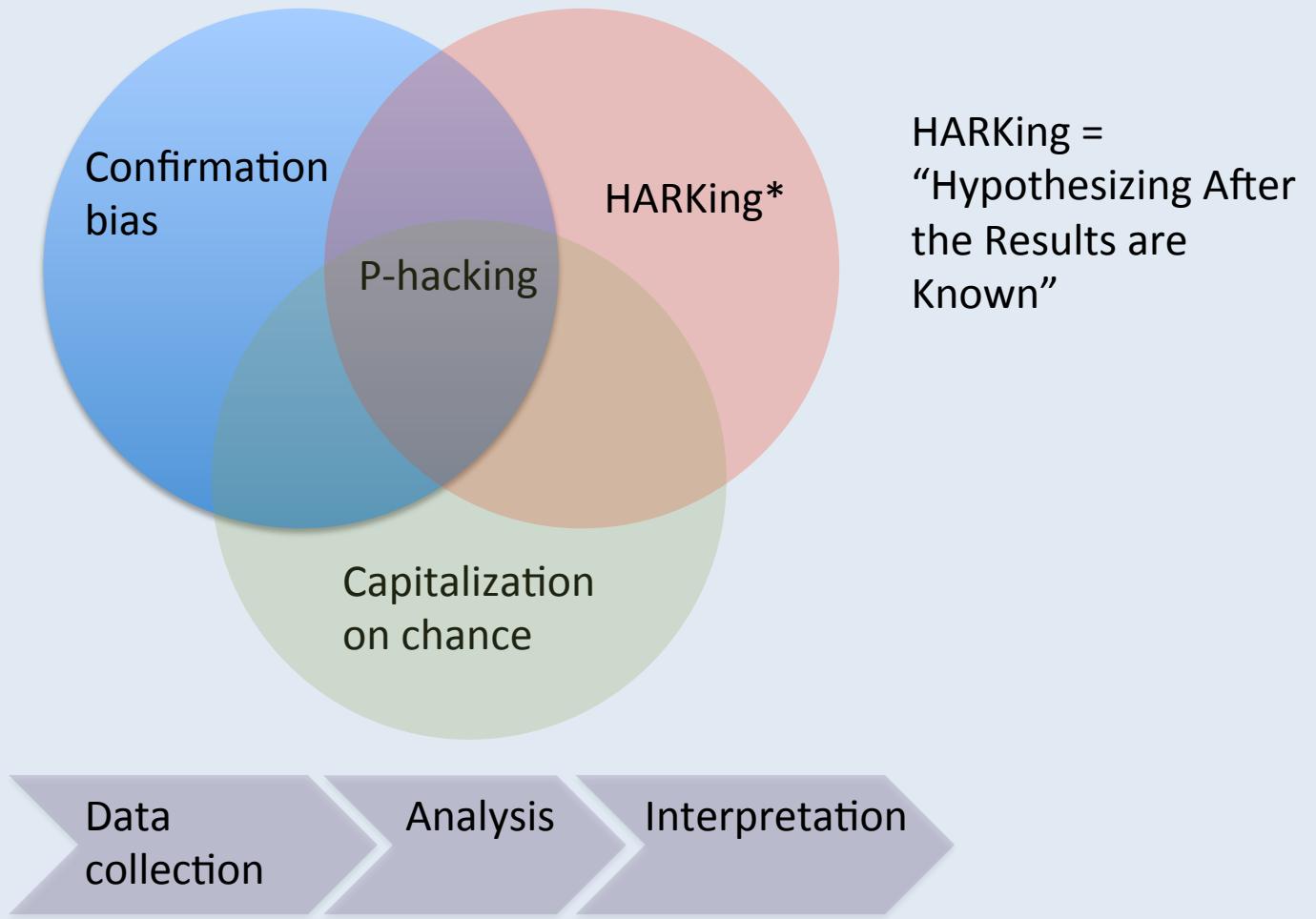
- *Biased assimilation:* Favored whichever study supported one's own views
- *Attitude polarization:* Groups grew farther apart after reading research

Disconfirmation Bias

- “But in *our* lab, we constantly seek out errors and artifacts, and we correct them.”
- Surely true,
 - But we are more likely to correct the results we didn't expect...
 - ...while overlooking errors that might be producing data that looks like what we expected.

“p hacking” to achieve statistical significance

- *Trimming outliers*
- *Collecting more cases*
 - *Throwing out cases*
 - *Subgroup analyses*



HARKing =
“Hypothesizing After
the Results are
Known”

SEXUAL ORIENTATION AND U.S. MILITARY PERSONNEL POLICY

An Update of RAND's 1993 Study

Drug War Heresies



Learning from Other
Vices, Times, & Places

Robert J. MacCoun
and Peter Reuter

CAMBRIDGE

more information · www.cambridge.org/0521572630

3/14/17



Report of the
Comprehensive Review
of the Issues Associated
with a Repeal of
“Don’t Ask, Don’t Tell”



November 30, 2010

Interpreting Dutch Cannabis Policy: Reasoning by Analogy in the Legalization Debate

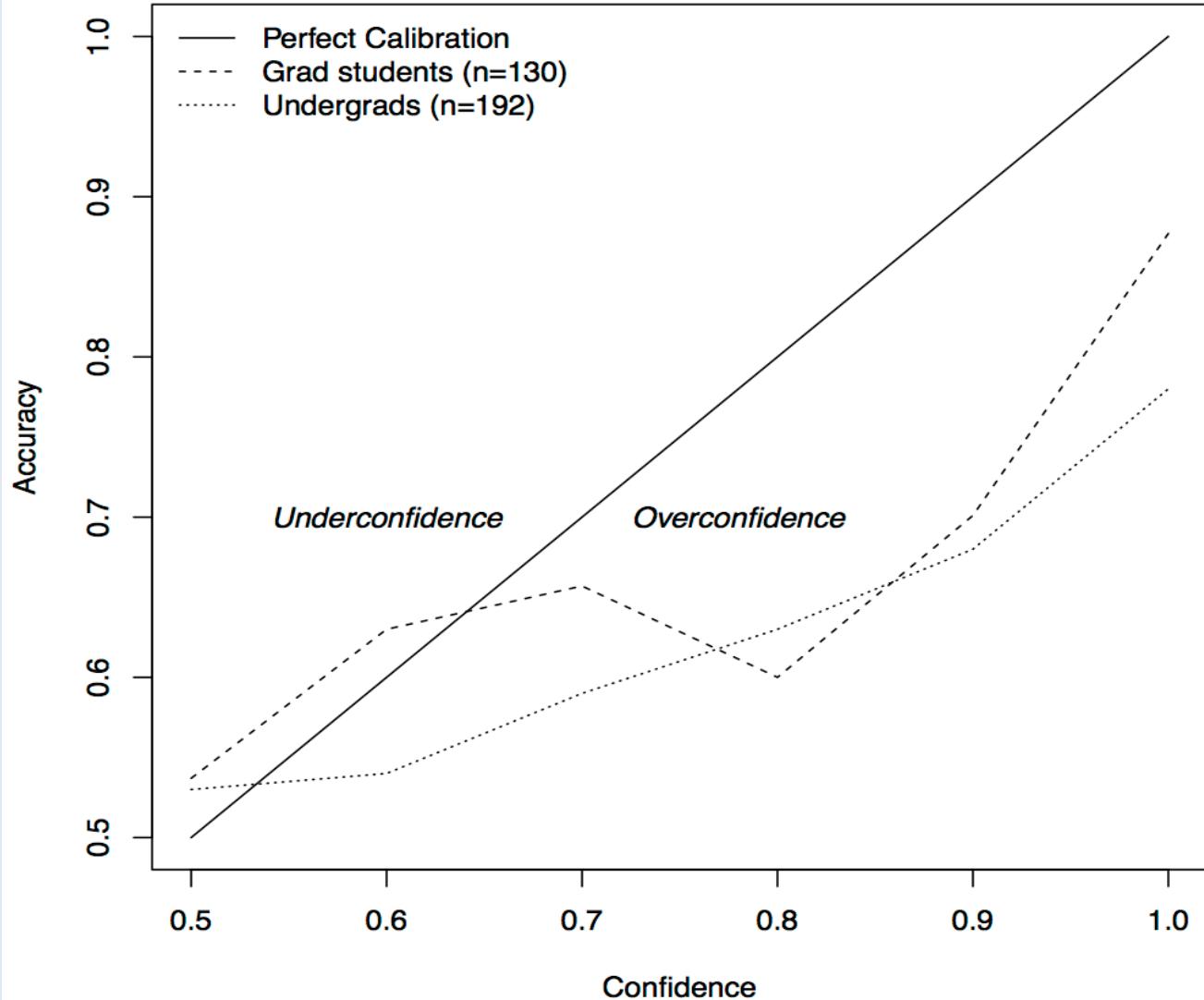
Robert MacCoun and Peter Reuter

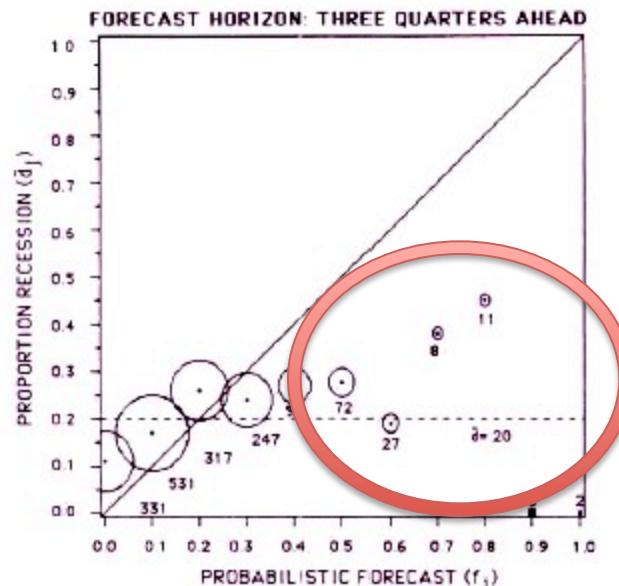
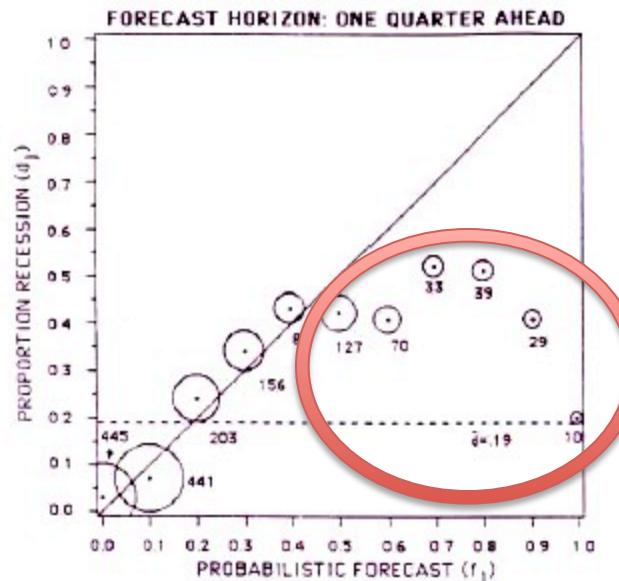
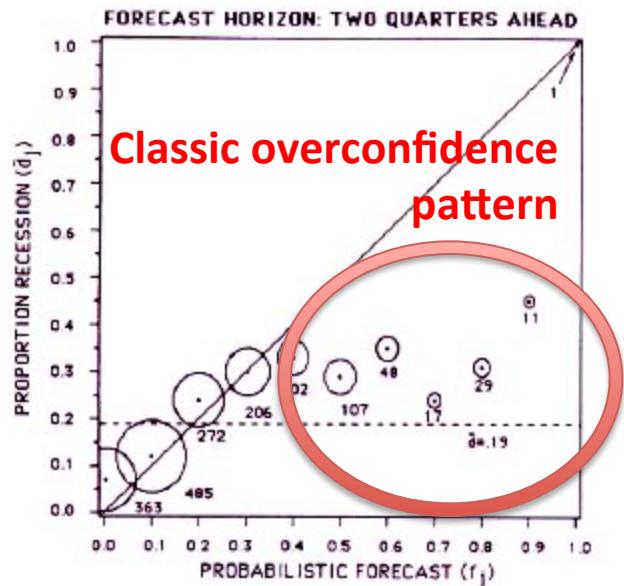
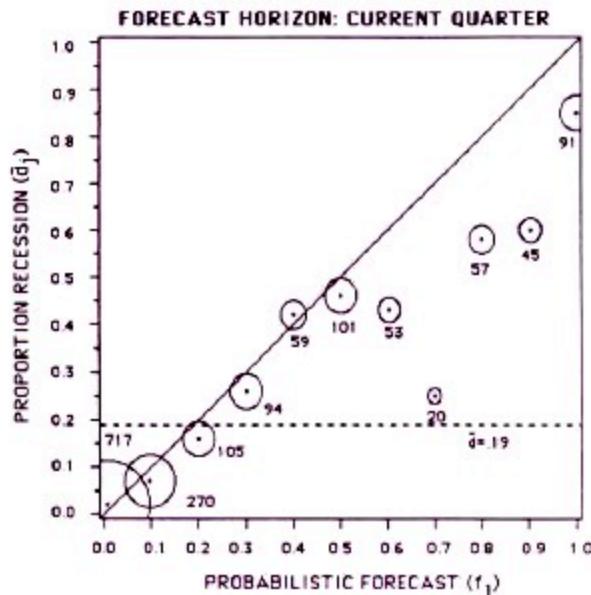
MacCoun

18

Expert overconfidence (Tetlock, 1998)

- Predictions made by professional foreign policy experts were “only slightly more accurate than one would expect from chance.”
- Average confidence rating across domains (where 9 = maximum confidence):
 - *Experts who were right:* 6.5 to 7.6
 - *Experts who were wrong:* 6.3 to 7.1





Braun & Yaniv, 1992

Survey panel of **professional economic forecasters** conducted quarterly by NBER and the Am. Statistical Assoc.

To play "Jeopardy!", Watson...builds an evidence profile to determine what are the most likely correct answers. Each answer has a confidence level.



For some questions, one answer will have a high confidence level. This is when Watson is most likely to buzz in. **For other questions, none of the answers will have a high confidence level and Watson will not buzz in.**

“That’s the interesting thing,” McQueeney said. “The machine knows when it doesn’t know the answer.”

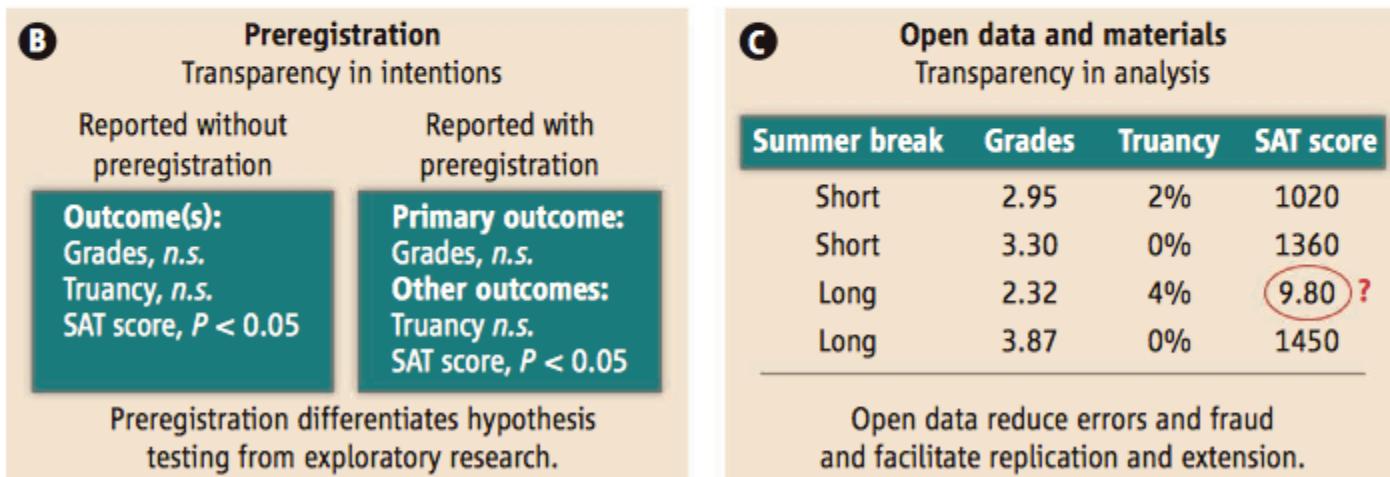
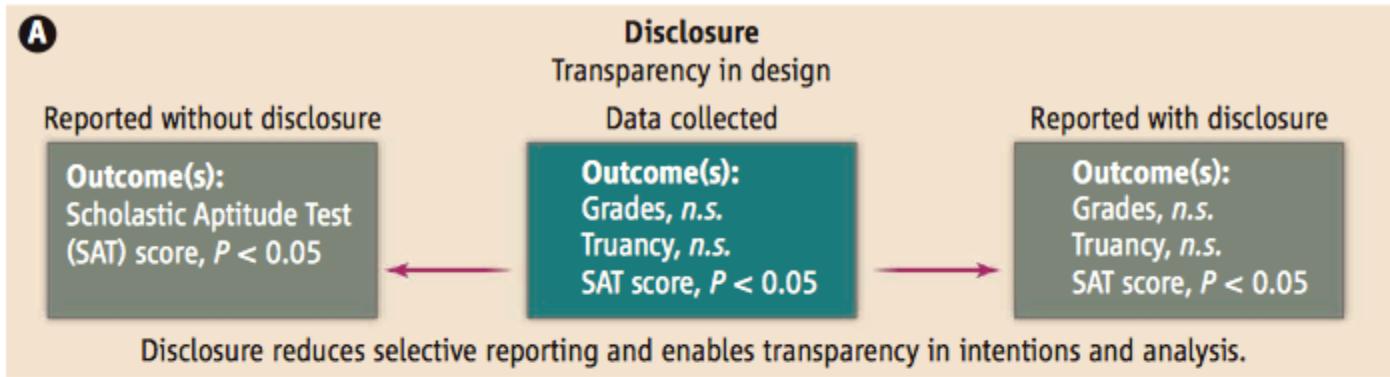
Corrective Practices

- Traditional practices
 - *peer reviewing*
 - *occasional “hired gun” replication*
 - *expert panels*
- New practices
 - *“Open science”*
 - *Preregistration of hypotheses*
 - *Routine “many labs” replication*
 - *Blind analysis*

Promoting Transparency in Social Science Research

3 JANUARY 2014 VOL 343 SCIENCE

E. Miguel,^{1*} C. Camerer,² K. Casey,³ J. Cohen,³ K. M. Esterling,⁴ A. Gerber,⁵ R. Glennerster,⁶ D. P. Green,⁷ M. Humphreys,⁷ G. Imbens,³ D. Laitin,³ T. Madon,¹ L. Nelson,¹ B. A. Nosek,^{8,9} M. Petersen,¹ R. Sedlmaier,¹⁰ J. P. Simmons,¹¹ U. Simonsohn,¹¹ M. Van der Laan¹



Three mechanisms for increasing transparency in scientific reporting. Demonstrated with a research question: "Do shorter summer breaks improve educational outcomes?" n.s. denotes $P > 0.05$.



Hide results to seek the truth

More fields should, like particle physics, adopt blind analysis to thwart bias, urge **Robert MacCoun** and **Saul Perlmutter**.

BLINDING STRATEGIES

Technique examples	Perturbation	Potential application
Noising $\theta_{ij} = y_{ij} + n_{ij}$ or $\theta_{ij} = \beta_k + n_{ij}$	Add a random number (from an appropriate statistical distribution) to data points or model parameters.	Testing which of several prevention messages is most effective in reducing smoking.
Biasing $\theta_{ij} = y_{ij} + b_j$	Obscure differences in experimental conditions by adding a hidden value that is biased in a particular direction.	Estimating whether the costs of a controversial safety regulation exceed its benefits.
Cell scrambling $\theta_{ij} = y_{\#}$	Shuffle labels for experimental conditions, so that it is unclear which set of results matches which conditions.	Testing a prediction that hard-copy books are better comprehended than audiobooks.
Item scrambling $\theta_{ij} = y_{##}$	Randomly relabel each data point to de-identify experimental conditions.	Analysing group differences that might be easy to recognize even with noise and bias (for example, effects of neighbourhood and school on crime victimization).
Various combinations	Row scrambling: keep pairs of variables together to preserve correlation. Variable blinding: swap labels of various variables.	

y_{ij} is the i th observation in the j th condition ('cell') of the study; β_k is the k th parameter of a model; θ_{ij} is y_{ij} or β_k after blinding; n_{ij} is random error, b_j is a bias term, and # denotes a randomly swapped subscript.

Different methods blind different features of the data

% of simulations where feature was undetectable:

	Add noise	Add bias	Scrambled cases	Scrambled cells
Range	100	100	0	0
Effect size	100	100	100	0
Exact p value	100	98	100	0
Stat. signif.	36	38	66	0
Direction	6	39	50	49

When should we blind?



When should we blind?



Key point to remember: The computer isn't blinded – you are.

The blinding can be limited to what's displayed rather than what's being analyzed.

When should we blind?



Or we can choose not to use blinding at all.

When should we blind?



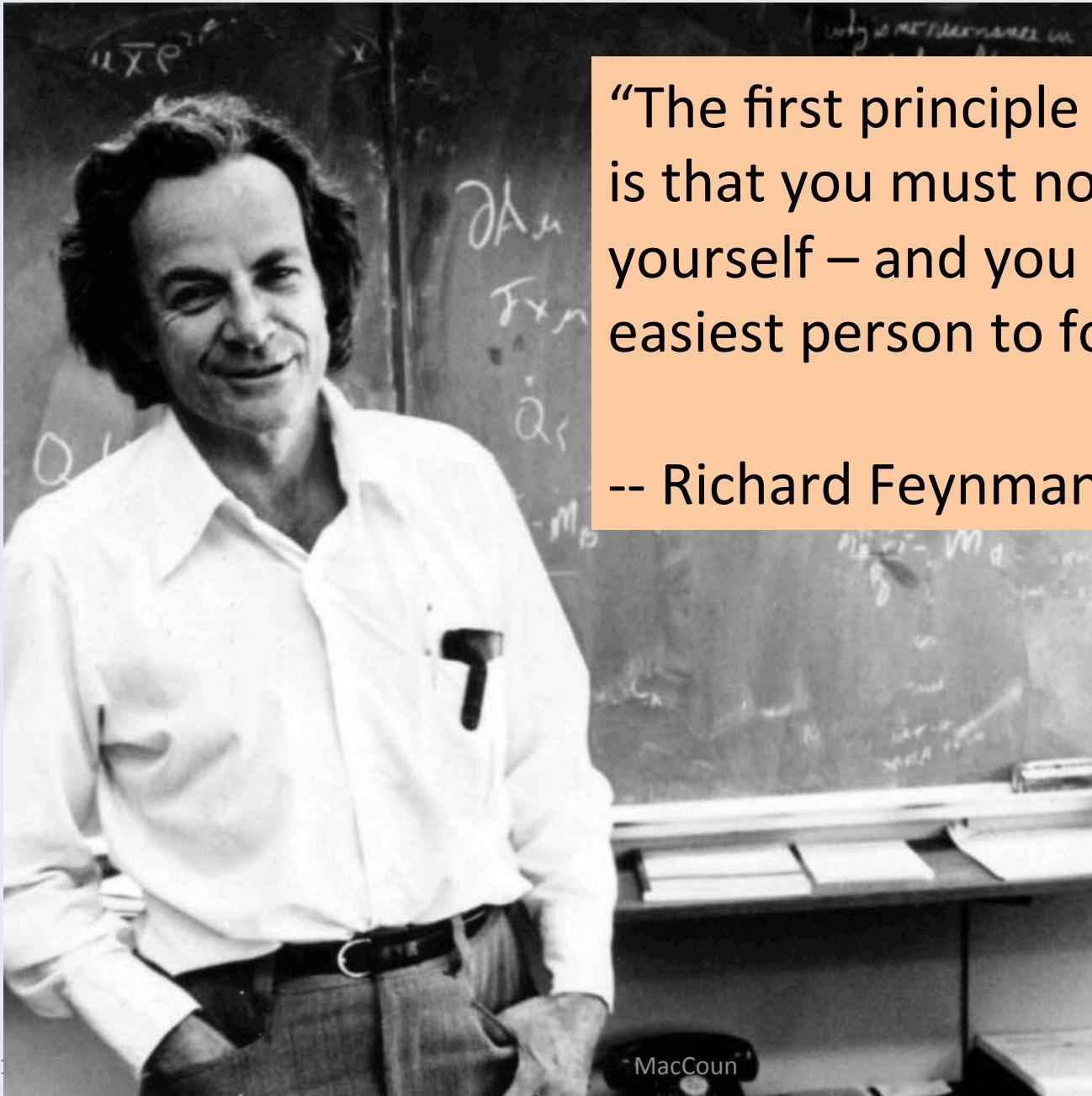
Or we can choose not to use blinding at all....

...but maybe we should start saying so:

“The data analysis was performed without blinding...”

Blind analysis should be tested empirically

- Benefits > costs?
- Comparison to other methods
 - Peer review, cross-validation, pre-registered analyses, cross-lab replication
- Does it impair discovery? Does it facilitate it?



“The first principle [of science] is that you must not fool yourself – and you are the easiest person to fool.”

-- Richard Feynman

A Poker Player Explains Why You Make Bad Business Decisions

What she just described is “outcome blind” analysis of decisions. Duke insists that all great poker players think this way. In fact she says she had a small peer group of elite players who regularly got together and talked shop. “When we told each other ‘the story of a hand,’ we wouldn’t mention whether we won or lost. It was all about the behaviors, the variables, the decisions – not the outcome.”



(Image courtesy of Annie Duke)

UNUSED SLIDES BELOW...

Canonical psychology experiment

- Advertisement urging people to vote for wetlands protection ballot initiative
- 2 x 2 factorial experiment (50 subjects per cell):
 - SOURCE EXPERTISE: 1=Low (source has a BA in Biology) vs. 2=High (source is a Yale Professor of Biology)
 - CONFLICT OF INTEREST: 3=None, 4=Yale will get new wetlands science center if initiative passes

To simulate raw data, we sampled from normal distribution:

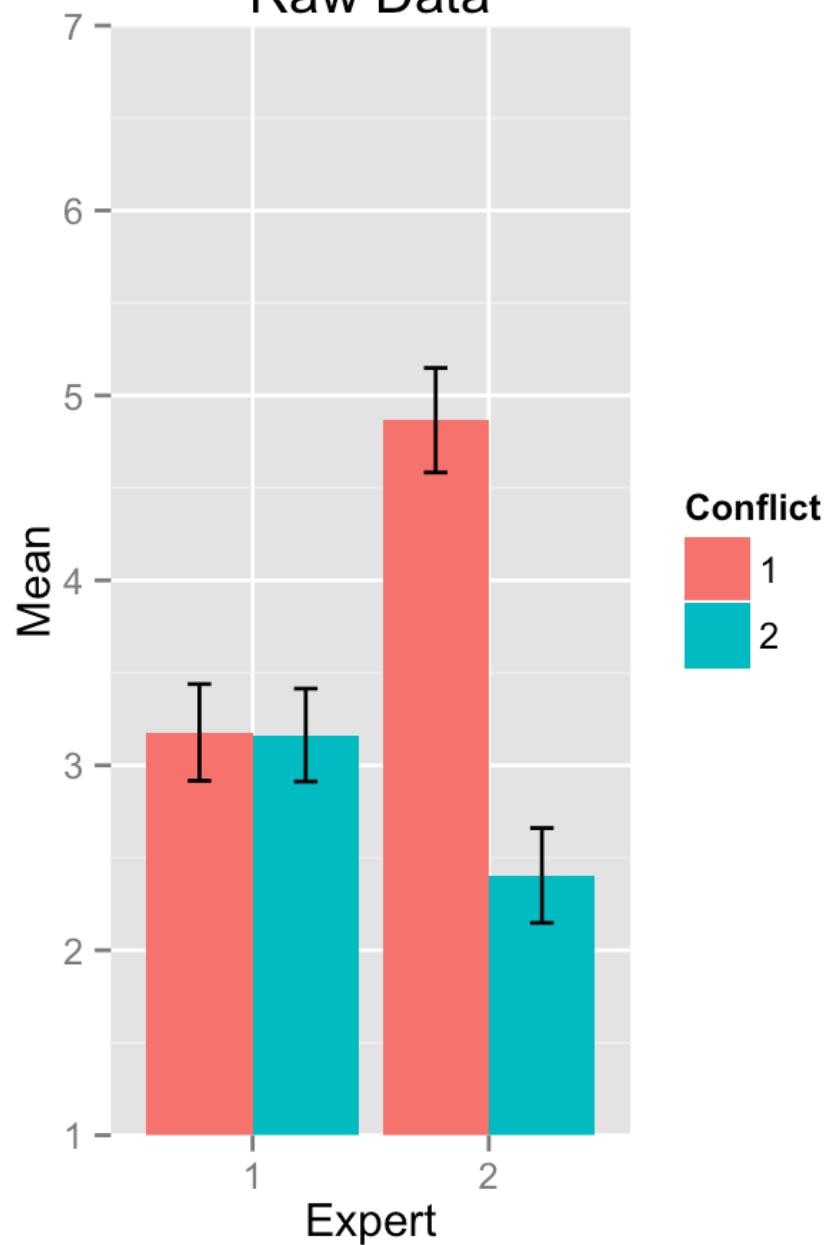
```
raw.DV00 <- rnorm(n=50,mean=3, sd=1) # Low Expert, No Conflict
```

```
raw.DV01 <- rnorm(n=50,mean=3, sd=1) # Low Expert, Conflict
```

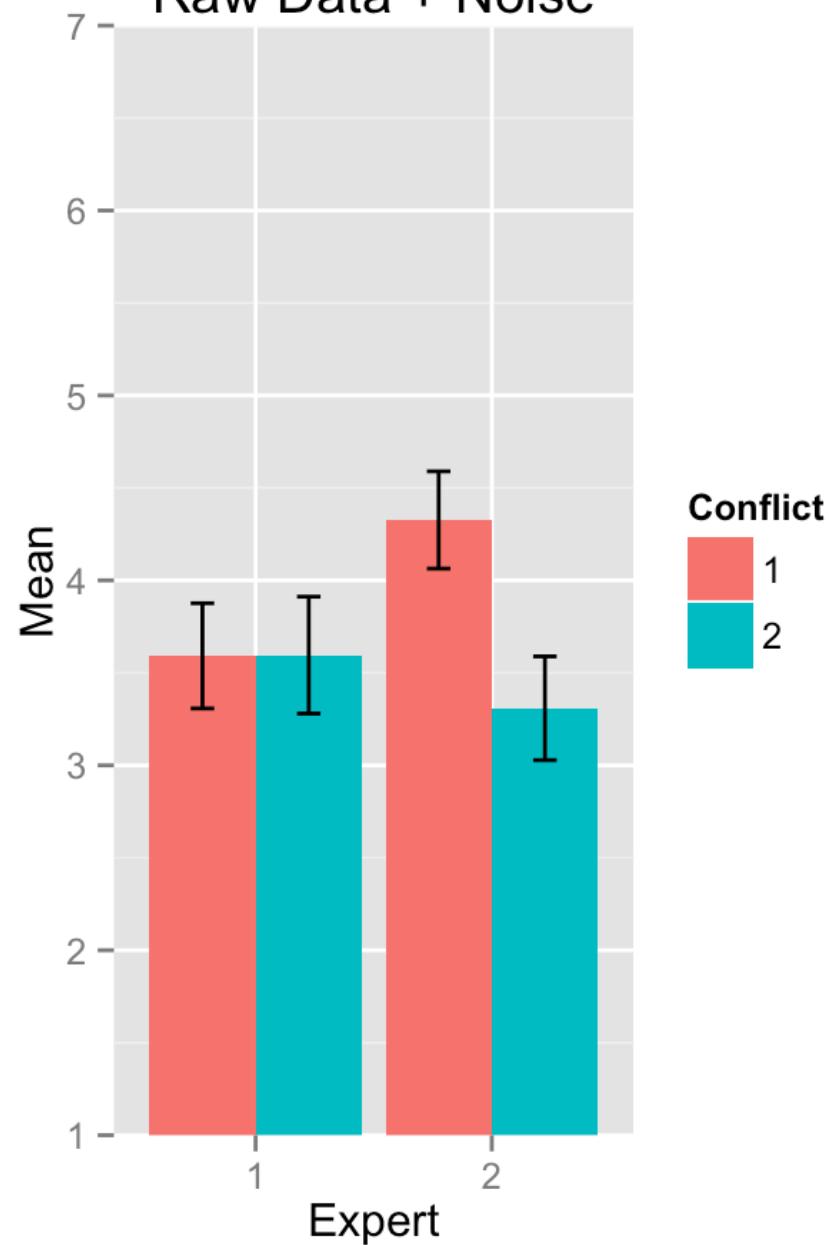
```
raw.DV10 <- rnorm(n=50,mean=4.5, sd=1) # High Expert, No Conflict
```

```
raw.DV11 <- rnorm(n=50,mean=2.5, sd=1) # High Expert, Conflict
```

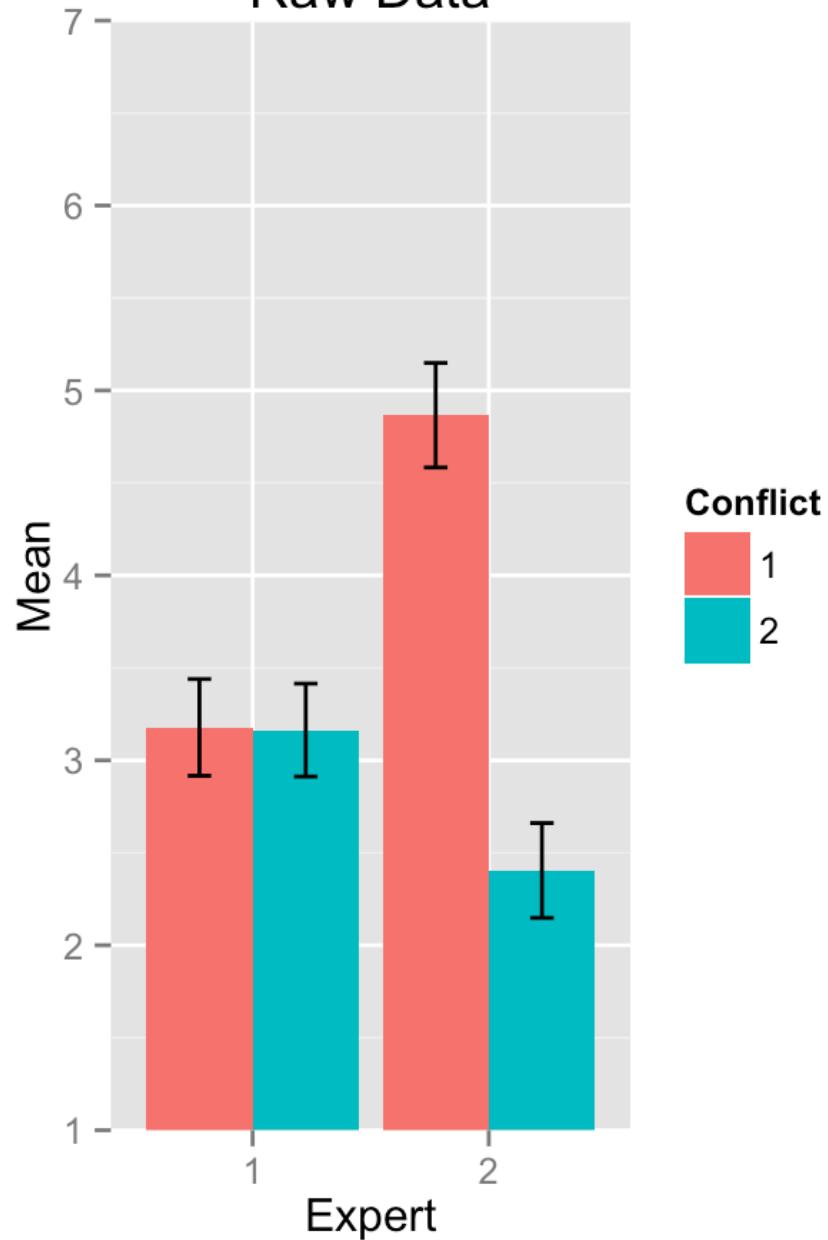
Raw Data



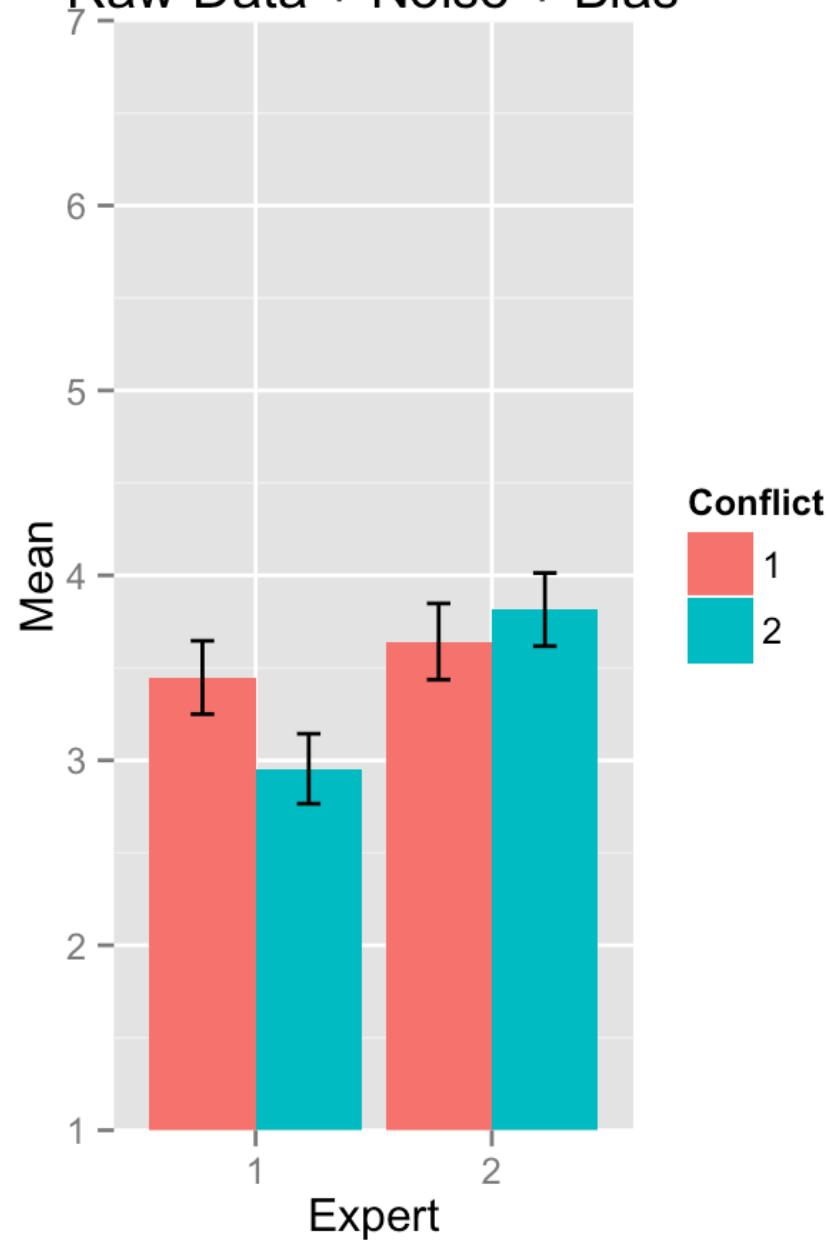
Raw Data + Noise



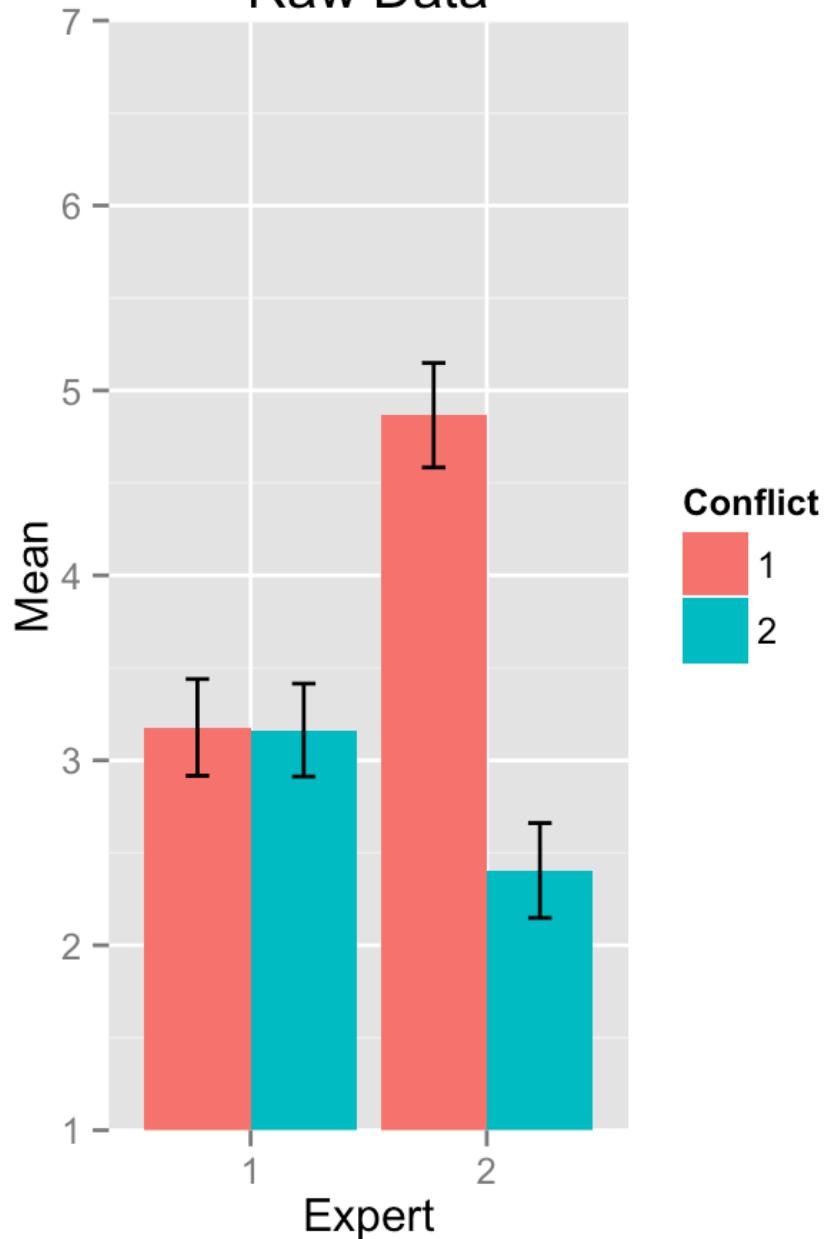
Raw Data



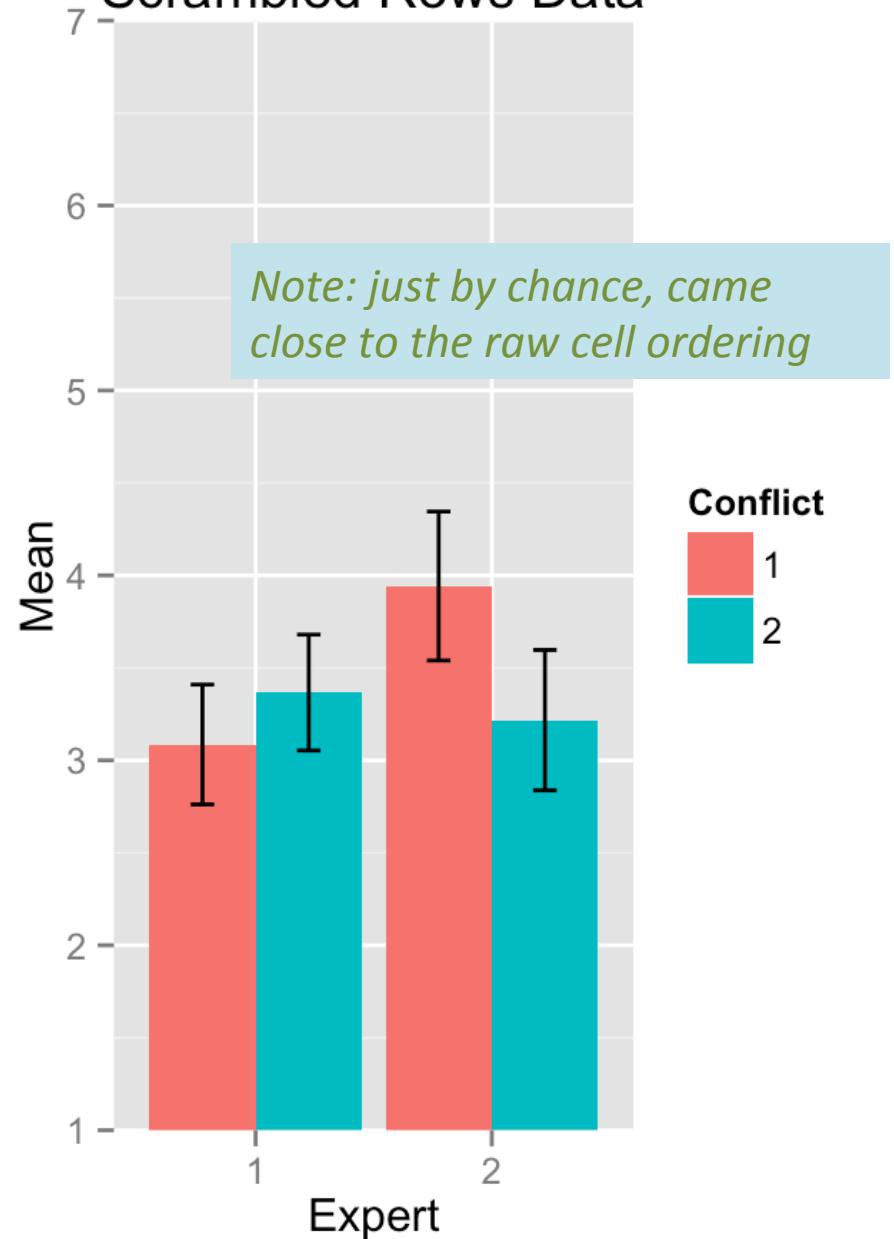
Raw Data + Noise + Bias



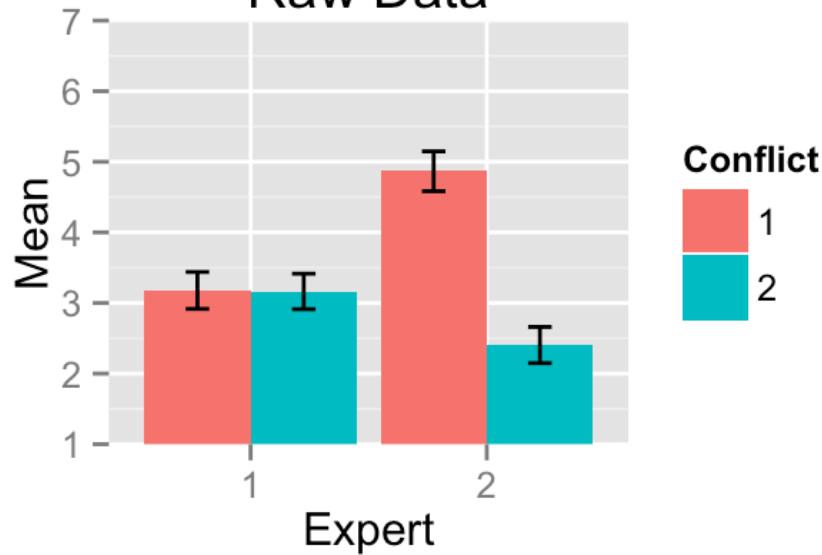
Raw Data



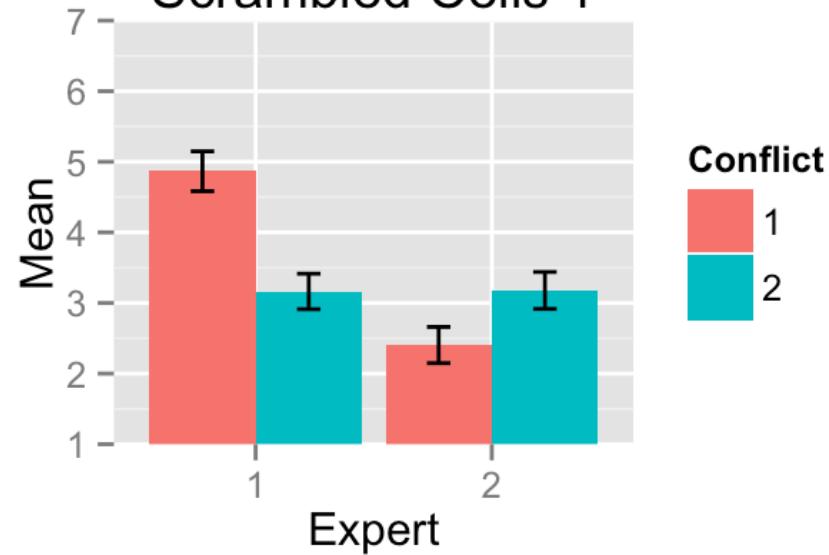
Scrambled Rows Data



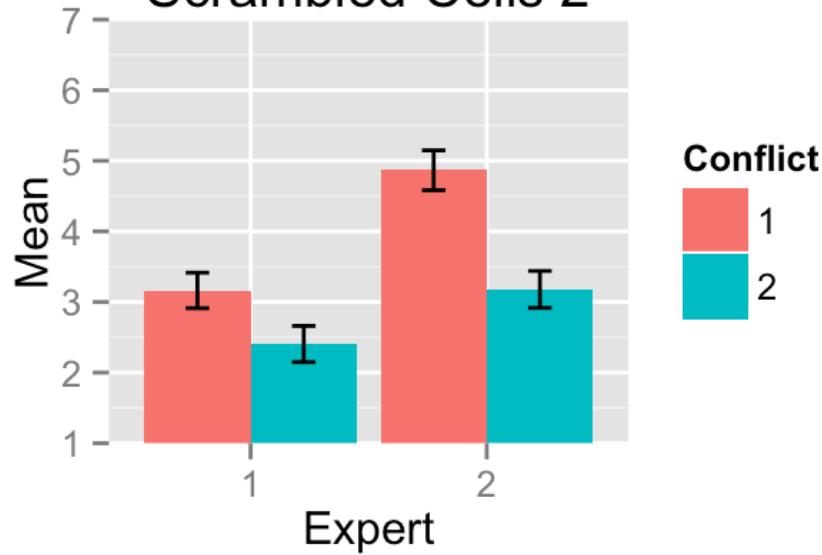
Raw Data



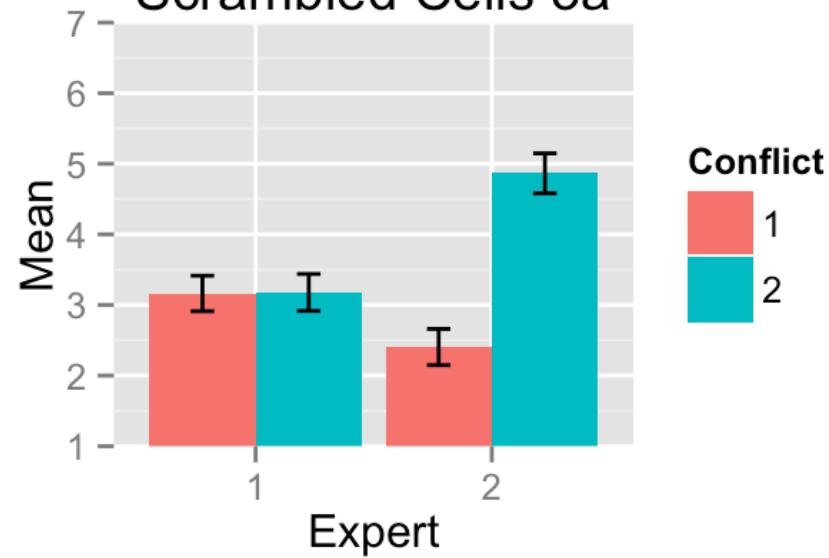
Scrambled Cells 1



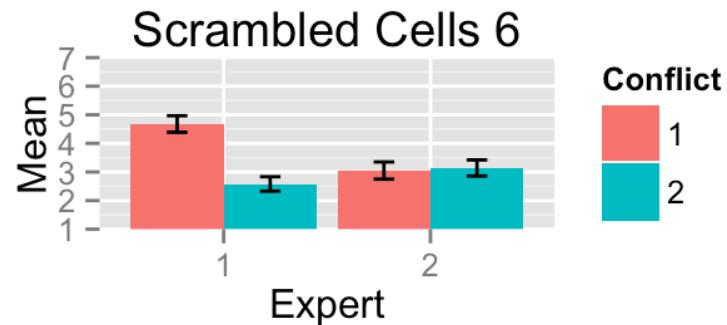
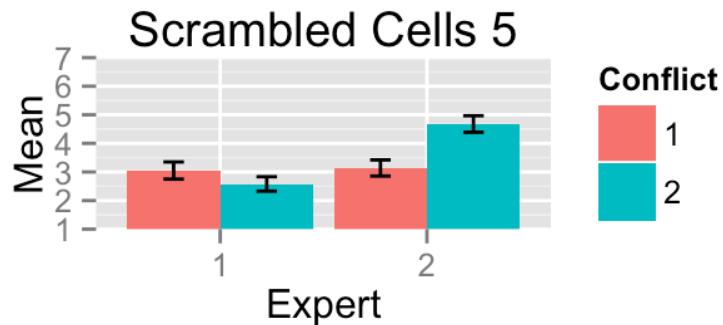
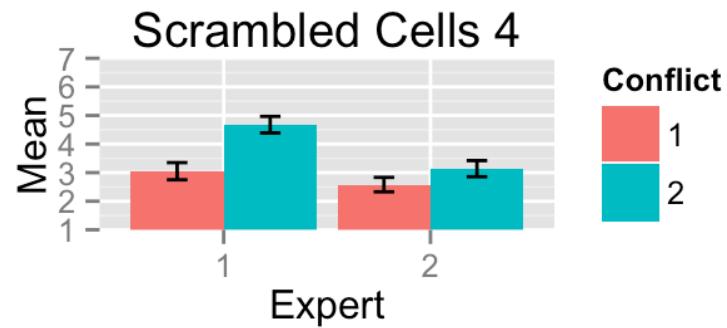
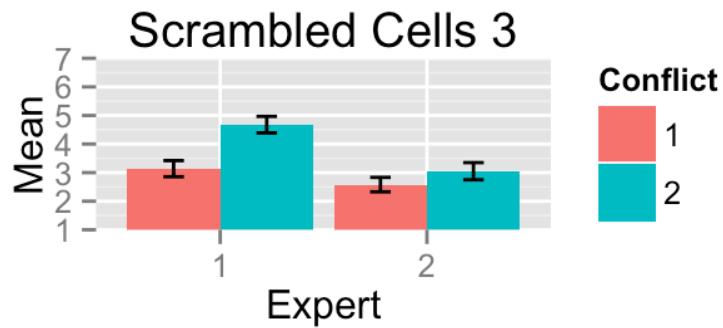
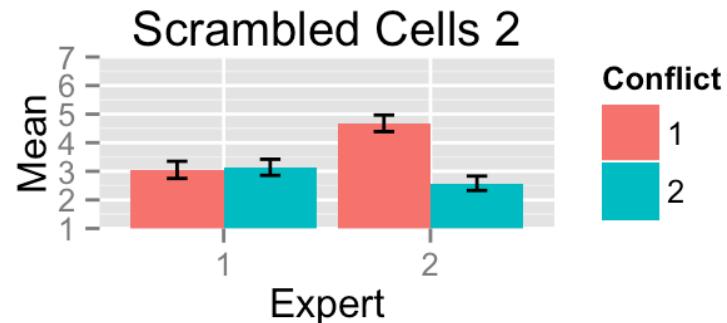
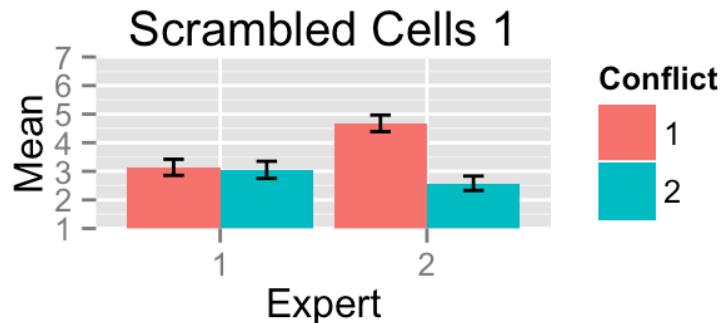
Scrambled Cells 2



Scrambled Cells 3a



Panel of 6 outcome sets, one of which is true....



Differences between physics and psychology research

- Physics has far higher degree of precision
- Physics has strong theory predicting quantities, and/or theoretical implications that can differ dramatically as a function of measured quantities
- Physics tends to have large research teams
- Physics data are sometimes sparse, difficult to obtain
- Random assignment to condition often impossible

2 x 2 ANOVA

	df	SS	MS	F	p
Expert	1	10.8	10.8	12.6	<.001
Conflict	1	76.6	76.6	89.3	<.001
E*C	1	74.9	74.9	87.3	<.001
Residual	196	168.1	0.9		

Blinding method 1: Add noise

- 200 random numbers sampled from uniform (min=1,max=7)
- $\text{blind.DV} = \text{mean}(\text{raw.DV}, \text{randnum})$

Raw data

	df	SS	MS	F	p
Expert	1	10.8	10.8	12.6	<.001
Conflict	1	76.6	76.6	89.3	<.001
E*C	1	74.9	74.9	87.3	<.001
Residual	196	168.1	0.9		

Blind data (raw + noise)

	df	SS	MS	F	p
Expert	1	2.5	2.5	2.5	0.119
Conflict	1	12.9	12.9	12.6	<.001
E*C	1	13.1	13.1	12.8	<.001
Residual	196	199.4	1.0		

Blinding method 2: Add noise + bias

- Same vector of random numbers as method 1
- Average together with a cell-specific bias term
 - Four bias terms, each sampled from $\text{normal}(\text{mean}=3, \text{sd}=1)$

Raw data

	df	SS	MS	F	p
Expert	1	10.8	10.8	12.6	<.001
Conflict	1	76.6	76.6	89.3	<.001
E*C	1	74.9	74.9	87.3	<.001
Residual	196	168.1	0.9		

Blind data (raw + noise + bias)

	df	SS	MS	F	p
Expert	1	13.9	13.9	28.7	<.001
Conflict	1	1.3	1.3	2.7	0.105
E*C	1	5.6	5.6	11.5	<.001
Residual	196	95.1	0.5		

Blinding method 3: Row scrambling

- There are $n=200$ rows in the data matrix
- “Row scrambling” simply randomly sorts their dependent variables, keeping the independent variables the same

Raw data

	df	SS	MS	F	p
Expert	1	10.8	10.8	12.6	<.001
Conflict	1	76.6	76.6	89.3	<.001
E*C	1	74.9	74.9	87.3	<.001
Residual	196	168.1	0.9		

Blind data (row scrambling)

	df	SS	MS	F	p
Expert	1	6.3	6.3	4.0	<.05
Conflict	1	2.5	2.5	1.6	0.21
E*C	1	12.7	12.7	8.0	<.01
Residual	196	309.1	1.6		

Comment on first three methods

- All three are *regressive* in the sense that the differences between conditions get weakened
- Problem:
 - *If data are noisy enough, no different than pre-registered data analysis (in the absence of any data)*
 - *Retaining more of the actual differences stimulates more thinking about the data*
 - *How to stimulate creativity without stimulating confirmatory bias?*

Blinding method 4: Cell scrambling

- Rather than scrambling individual data points, this method keeps cell data intact, but scrambles the four cells of the design
- In this case, there are $4! = 24$ possible orderings, and we've randomly sampled three of them...
- In all 3 cases, the main effects and interaction are significant ($p < .001$), but the interpretation differs...

Analytic Method	What Could Be Blinded?
Analysis of variance for randomized controlled trials	<ul style="list-style-type: none"> • Outcome values • Statistical significance • Experimental condition
Ordinary least-square regression analysis	<ul style="list-style-type: none"> • Outcome values • Statistical significance • Variable labels
Factor analysis	<ul style="list-style-type: none"> • Outcome values • Statistical significance • Variable labels
Time series analysis	<ul style="list-style-type: none"> • Outcome values • Statistical significance • Event dates

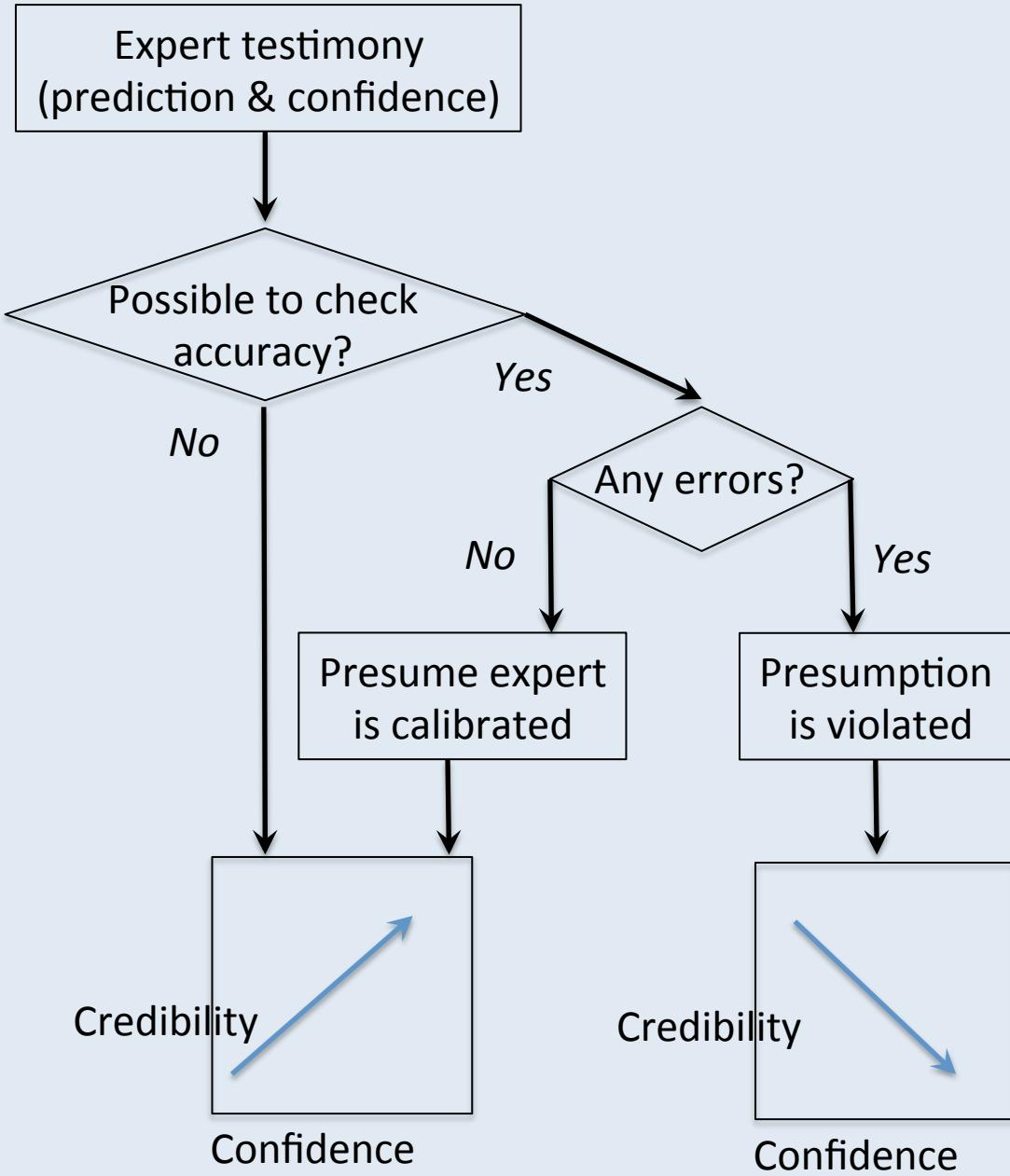


Table 1. Citizens' motives and their perceptions of expert motives.

		Perceived expert motive	
		<i>Inquisitorial</i>	<i>Adversarial</i>
Citizen's motives	<i>Inquisitorial</i>	Who is more likely to be correct?	Who is more honest?
	<i>Adversarial</i>	Does this help or hurt our side?	



Political Psychology, Vol. 30, No. 1, 2009

Citizens' Perceptions of Ideological Bias in Research on Public Policy Controversies

Robert J. MacCoun

University of California, Berkeley

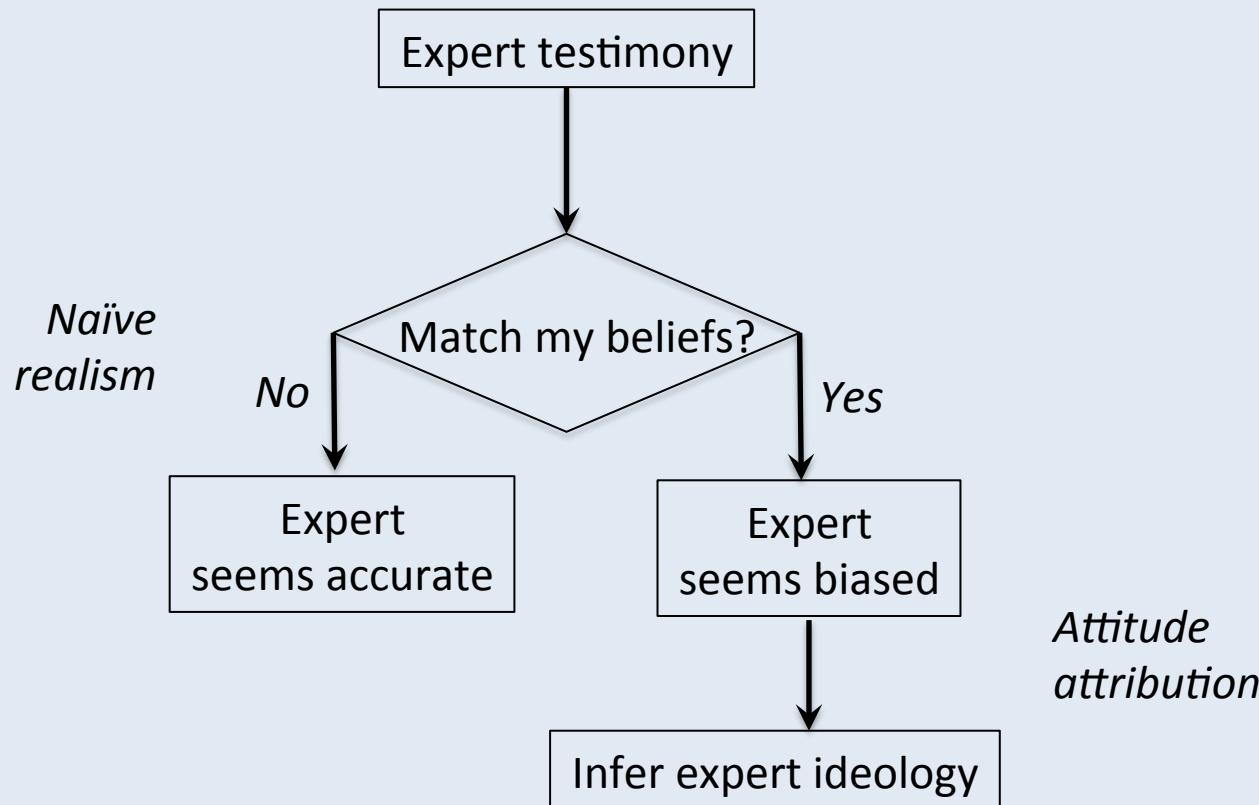
Susannah Paletz

University of California, Berkeley

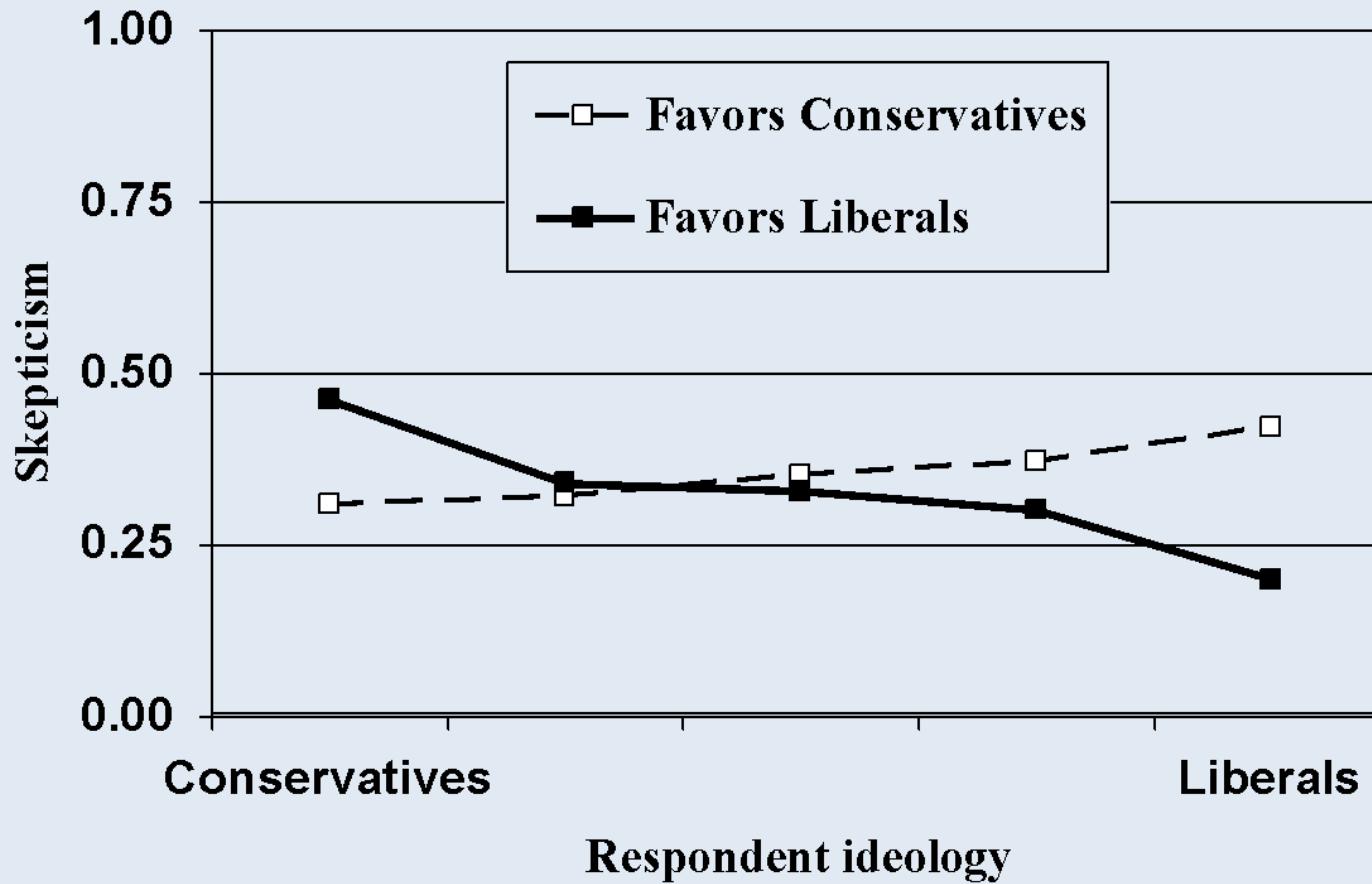


Table 1. Research Design and Cell Sizes

	Finding favors conservative position	Finding favors liberal position
Gun control	Fails to reduce crime (n = 120)	Reduces crime (n = 104)
Death penalty	Deters crime (n = 90)	Fails to deter crime (n = 122)
Medical marijuana	No medical value (n = 93)	Has medical value (n = 88)
School vouchers	Improves education (n = 124)	Fails to improve education (n = 105)
Nutrition ads	TV Actor (n = 102) vs. Real Doctor (n = 102)	



How doubtful is result?



How liberal is the researcher?

