

# Preregister Experiments with AI Agents

**Michelle Vaccaro**

MIT Institute for Data, Systems, and Society

# People are treating AI agents as participants in experiments

## Using Large Language Models to Simulate Multiple Humans and Replicate Human Subject Studies

Gati Aher<sup>1</sup> Rosa I. Arriaga<sup>2</sup> Adam Tauman Kalai<sup>3</sup>



## Evaluating and Inducing Personality in Pre-trained Language Models

Guangyuan Jiang<sup>1,2,\*</sup>  
jgy@stu.pku.edu.cn

Manjie Xu<sup>1,\*</sup>  
manjietsu@gmail.com

Song-Chun Zhu<sup>1,3</sup>  
s.c.zhu@pku.edu.cn

Wenjuan Han<sup>4,✉</sup>  
wjhan@bjtu.edu.cn

Chi Zhang<sup>3,✉</sup>  
zhangchi@bigai.ai

Yixin Zhu<sup>1,✉</sup>  
yixin.zhu@pku.edu.cn

## Evaluating the Moral Beliefs Encoded in LLMs

**Warning:** This paper contains moral scenarios which are controversial and offensive in nature.

Nino Scherrer<sup>\*1</sup>, Claudia Shi<sup>\*1,2</sup>, Amir Feder<sup>2</sup>, and David M. Blei<sup>2</sup>

<sup>1</sup> FAR AI, <sup>2</sup> Columbia University

## MACHINE PSYCHOLOGY

Thilo Hagendorff<sup>†</sup>  
University of Stuttgart

Ishita Dasgupta<sup>\*</sup>  
Google DeepMind

Marcel Binz<sup>†</sup>  
Helmholtz Institute for  
Human-Centered AI

Stephanie C.Y. Chan<sup>†</sup>  
Google DeepMind

Andrew Lampinen<sup>†</sup>  
Google DeepMind

Jane X. Wang<sup>†</sup>  
Google DeepMind

Zeynep Akata  
TU Munich

Eric Schulz  
Helmholtz Institute for  
Human-Centered AI

## EVALUATING LANGUAGE MODEL AGENCY THROUGH NEGOTIATIONS

Tim R. Davidson<sup>†\*</sup> Veniamin Veselovsky<sup>†\*</sup> Martin Josifoski<sup>†</sup> Maxime Peyrard<sup>‡</sup>

Antoine Bosselut<sup>†</sup> Michal Kosinski<sup>§</sup> Robert West<sup>†</sup>

<sup>†</sup>EPFL, <sup>‡</sup>UGA, CNRS, LIG, <sup>§</sup>Stanford University

## Instructed to Bias: Instruction-Tuned Language Models Exhibit Emergent Cognitive Bias

Itay Itzhak<sup>1</sup>, Gabriel Stanovsky<sup>2</sup>, Nir Rosenfeld<sup>1</sup>, Yonatan Belinkov<sup>1</sup>

<sup>1</sup>Technion – Israel Institute of Technology, Israel

<sup>2</sup>School of Computer Science and Engineering, The Hebrew University of Jerusalem

itaylitzhak@gmail.com,

{nirr, belinkov}@technion.ac.il, gabriel.stanovsky@mail.huji.ac.il

# Experiments with AI agents exacerbate vulnerabilities that exist in experiments with human subjects

> [Psychol Sci.](#) 2011 Nov;22(11):1359-66. doi: 10.1177/0956797611417632. Epub 2011 Oct 17.

## False-positive psychology: undisclosed flexibility in data collection and analysis allows presenting anything as significant

Joseph P Simmons <sup>1</sup>, Leif D Nelson, Uri Simonsohn

Affiliations + expand

PMID: 22006061 DOI: [10.1177/0956797611417632](#) 

### Abstract

In this article, we accomplish two things. First, we show that despite empirical psychologists' nominal endorsement of a low rate of false-positive findings ( $\leq .05$ ), flexibility in data collection, analysis, and reporting dramatically increases actual false-positive rates. In many cases, a researcher is more likely to falsely find evidence that an effect exists than to correctly find evidence that it does not. We present computer simulations and a pair of actual experiments that demonstrate how unacceptably easy it is to accumulate (and report) statistically significant evidence for a false hypothesis. Second, we suggest a simple, low-cost, and straightforwardly effective disclosure-based solution to this problem. The solution involves six concrete requirements for authors and four guidelines for reviewers, all of which impose a minimal burden on the publication process.

> [Psychol Sci.](#) 2012 May 1;23(5):524-32. doi: 10.1177/0956797611430953. Epub 2012 Apr 16.

## Measuring the prevalence of questionable research practices with incentives for truth telling

Leslie K John <sup>1</sup>, George Loewenstein, Drazen Prelec

Affiliations + expand

PMID: 22508865 DOI: [10.1177/0956797611430953](#) 

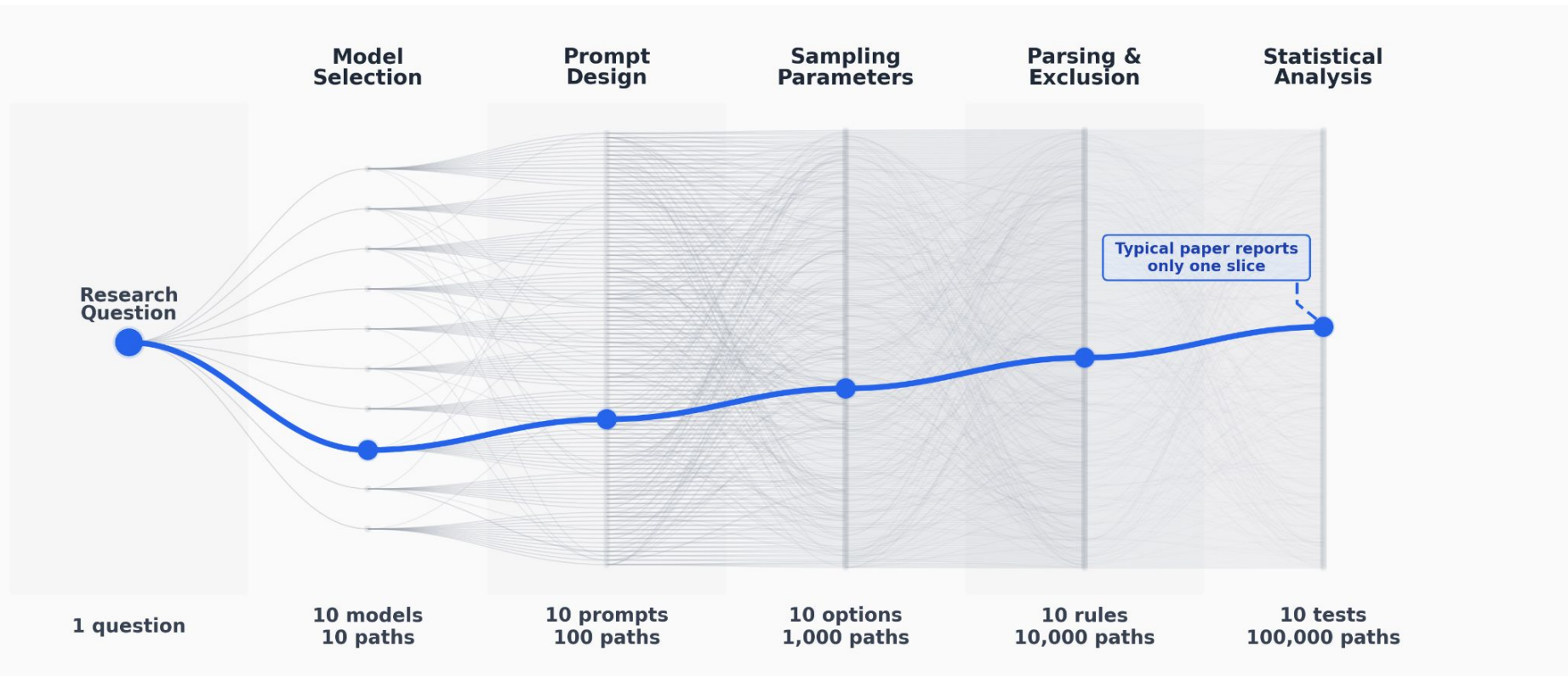
+ Paperpile

### Abstract

Cases of clear scientific misconduct have received significant media attention recently, but less flagrantly questionable research practices may be more prevalent and, ultimately, more damaging to the academic enterprise. Using an anonymous elicitation format supplemented by incentives for honest reporting, we surveyed over 2,000 psychologists about their involvement in questionable research practices. The impact of truth-telling incentives on self-admissions of questionable research practices was positive, and this impact was greater for practices that respondents judged to be less defensible. Combining three different estimation methods, we found that the percentage of respondents who have engaged in questionable practices was surprisingly high. This finding suggests that some questionable practices may constitute the prevailing research norm.

**HARKing, *p*-hacking, design-hacking, selective reporting...**

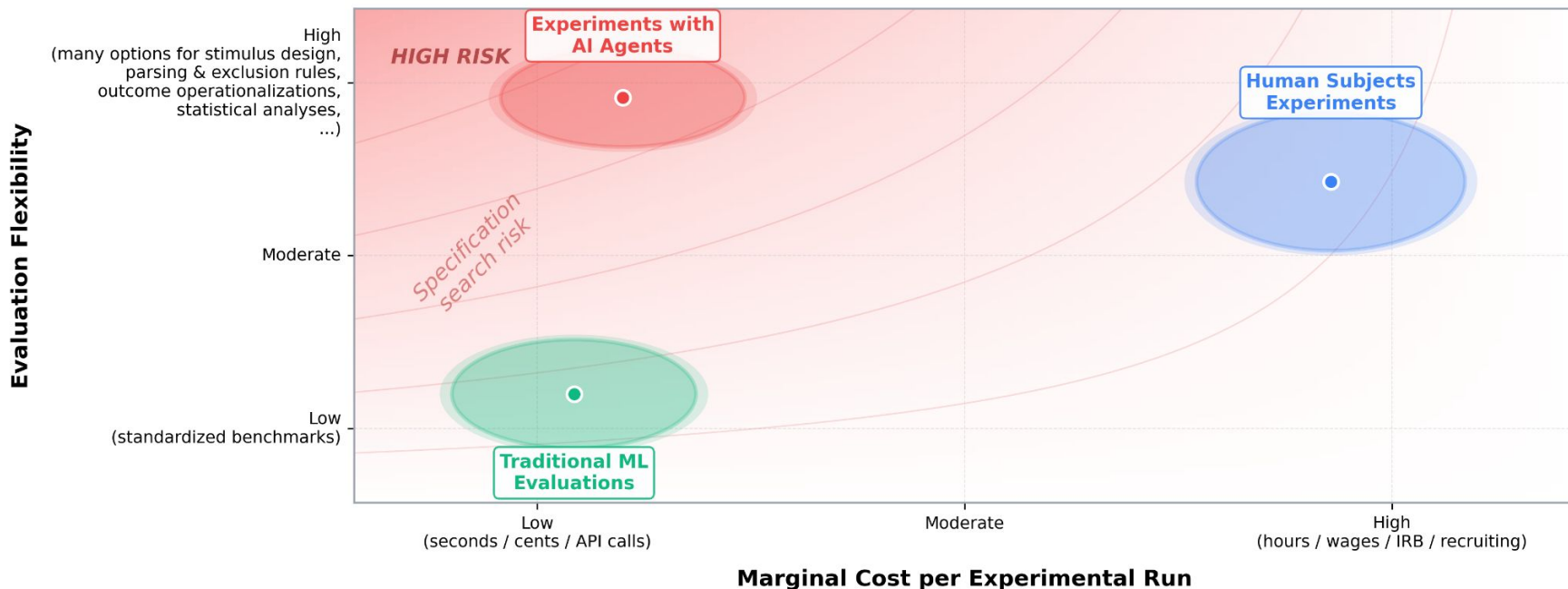
# Experiments with AI agents exacerbate vulnerabilities that exist in experiments with human subjects



**Combinatorial explosion of researcher degrees of freedom**

# Experiments with AI agents exacerbate vulnerabilities that exist in experiments with human subjects

## Cost-Flexibility Tradeoff Across Research Paradigms



**Low experimental cost, high evaluation flexibility = high specification search risk**

# We should preregister experiments with AI agents to address these issues

RESEARCH ARTICLE | PSYCHOLOGICAL AND COGNITIVE SCIENCES | 

## The preregistration revolution

[Brian A. Nosek](#)  , [Charles R. Ebersole](#) , [Alexander C. DeHaven](#) , and [David T. Mellor](#)  [Authors Info & Affiliations](#)

Edited by Richard M. Shiffrin, Indiana University, Bloomington, IN, and approved August 28, 2017 (received for review June 15, 2017)

March 12, 2018 | 115 (11) 2600-2606 | <https://doi.org/10.1073/pnas.1708274114>

## Preregistration Template for Experiments with AI Agents

Version 1.0

### Purpose

This template is designed to facilitate preregistration of experiments that use AI agents (e.g., large language models) as participants. It addresses the unique researcher degrees of freedom in AI experimentation: model selection, prompt engineering, sampling parameters, response processing, and analysis choices. Complete this template before data collection begins and submit to a public registry (e.g., OSF Registries).

**Extending the logic of traditional preregistration for the new degrees of freedom that characterize AI agent research**

# Thank you!

**Questions? Comments? Suggestions?**  
Please leave a comment or reach out to me!