

1 A survey of registration practices among observational researchers using preexisting datasets

2 Robert T. Thibault<sup>1,5</sup>, Marton Kovacs<sup>2,6</sup>, Tom E. Hardwicke<sup>3</sup>, Alexandra Sarafoglou<sup>4</sup>, John  
3 P. A. Ioannidis<sup>4</sup>, & Marcus R. Munafò<sup>1,7</sup>

4 <sup>1</sup> Meta-Research Innovation Center at Stanford (METRICS), Stanford University.

5 <sup>2</sup> Doctoral School of Psychology, ELTE Eotvos Lorand University, Budapest, Hungary

6 <sup>3</sup> Melbourne School of Psychological Sciences, University of Melbourne.

7 <sup>4</sup> Department of Psychology, University of Amsterdam.

8 <sup>5</sup> School of Psychological Science, University of Bristol.

9 <sup>6</sup> Institute of Psychology, ELTE Eotvos Lorand University, Budapest, Hungary

10 <sup>7</sup> Meta-Research Innovation Center Berlin (METRIC-B), QUEST Center for Transforming  
11 Biomedical Research, Berlin Institute of Health, Charité – Universitätsmedizin Berlin.

12 <sup>8</sup> MRC Integrative Epidemiology Unit at the University of Bristol.

13 <sup>9</sup> Departments of Medicine, Epidemiology and Population Health, Biomedical Data Science,  
14 and Statistics, Stanford University.

The authors made the following contributions. Robert T. Thibault: Conceptualization, Data curation, Formal analysis, Funding acquisition, Investigation, Methodology, Project administration, Resources, Supervision, Validation, Visualization, Writing - original draft, Writing - review & editing; Marton Kovacs: Data curation, Formal analysis, Software, Validation, Visualization, Writing - review & editing; Tom E. Hardwicke: Methodology, Writing - review & editing; Alexandra Sarafoglou: Methodology, Writing - review & editing; John P. A. Ioannidis: Methodology, Writing - review & editing; Marcus R. Munafò: Conceptualization, Methodology, Supervision, Writing - review & editing.

Correspondence concerning this article should be addressed to Robert T. Thibault, Enter postal address here. E-mail: robert.thibault@stanford.edu

26

Abstract

27 placeholder for an abstract

28 *Keywords:* keywords

29 Word count: X

30 A survey of registration practices among observational researchers using preexisting datasets

## 31 Introduction

## 32 Methods

## 33 Results

### 34 Participants

35 We invited the ALSPAC mailing list to participate, which included 1148 email  
36 addresses. 54 emails bounced, leaving 1094 emails that went through. The survey was  
37 completed 103 times and partially completed 20 times, leading to a response rate of 11% for  
38 complete surveys and 2% for incomplete surveys.<sup>1</sup> The median time taken for complete  
39 survey responses was 7.40 minutes (IQR: 4.60 to 13.10).

40 Respondents published a median of NA (IQR 2 to 26.20) studies using preexisting  
41 observational data (Figure S1). They reported using the programming languages R (n = 65),  
42 Stata (n = 48), SPSS (n = 17), SAS (n = 15), Python (n = 6), Mplus (n = 3), Bash (n = 2),  
43 MATLAB (n = 1), Nextflow (n = 1), and plink2 (n = 1) (Table S1)<sup>2</sup>. 62% (62/100) of  
44 participants reported being more concerned with research trustworthiness, bias, rigour, and  
45 reproducibility compared to what they think of as a typical research who uses preexisting  
46 observational data (Figure C2); 6% (6/100) reported being less concerned.

---

<sup>1</sup> The ALSPAC mailing list has been active for >30 years and may contain email addresses that are no longer monitored. For example, we received one email reply stating that the recipient hasn't been active in research for 30 years. Excluding these email addresses would increase the response rate, but we do not know by how much.

<sup>2</sup> Participants could select multiple responses to this survey question.

## Survey results

Most respondents agreed that studies that analyze preexisting observational datasets are trustworthy<sup>3</sup> (72%; 74/103) and reproducible<sup>4</sup> (79%; 81/103) (Figure 2, top panel). At the same time, many agreed that a study using an ECAW would be *more* trustworthy (70%; 70/100) and *more* reproducible (68%; 69/101) compared to a typical study using preexisting observational data (Figure 2, bottom panel).

Over half of respondents reported that their studies using preexisting observational data are preregistered never or almost never (36%; 37/103), or sometimes (25%; 26/103) (Figure 2A). About half reported sharing their analysis scripts never or almost never (20%; 21/103), or sometimes (32%; 33/103) (Figure 3). 77% (79/103) reported that they never or almost never blind the data analyst (Figure 3). Almost all respondents answered that they use both exploratory (93%; 96/103) and confirmatory (87%; 90/103) analyses at least sometimes (Figure 3).

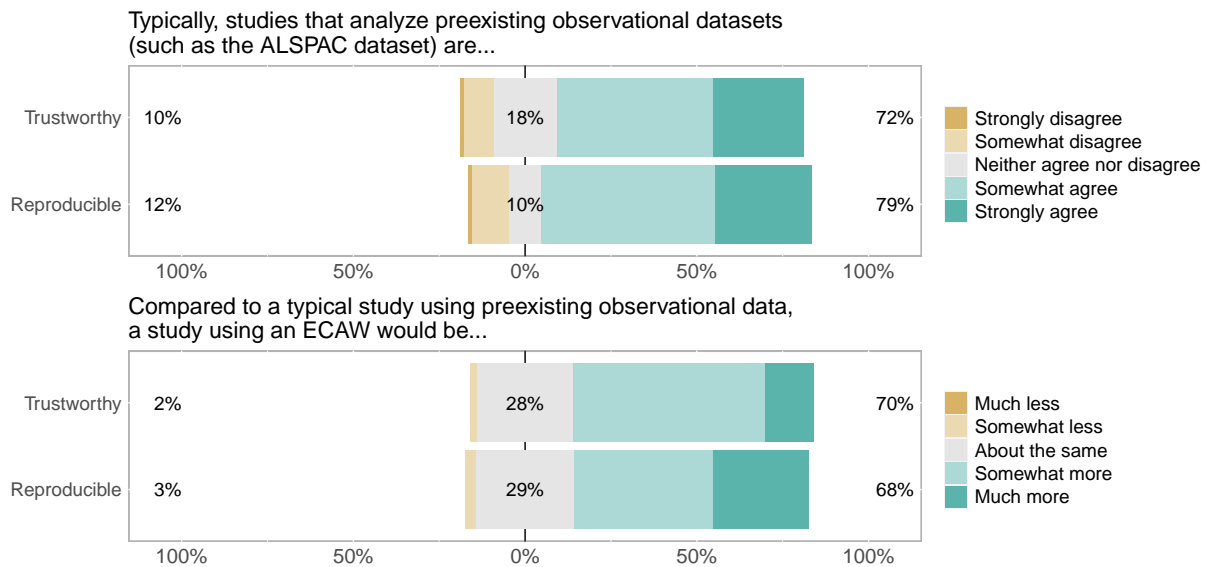
26% (26/101) of respondents agreed (versus 45%; 45/101 who disagreed) that they would be less willing to use ALSPAC data if they were required to use an ECAW (Figure 3). 53% (50/94) agreed (20%; 19/94 disagreed) that they would opt-in if ALSPAC ran a study on ECAWs. 55% (53/96) agreed (10%; 10/96 disagreed) that ALSPAC should run a study on ECAWs. 46% (43/94) agreed (22%; 21/94 disagreed) that they would prefer using an ECAW than using typical preregistration.

Table 1. Recurring topics in responses to the open-ended survey questions. The survey included 4 open-ended questions with broad prompts regarding running a study on ECAWs,

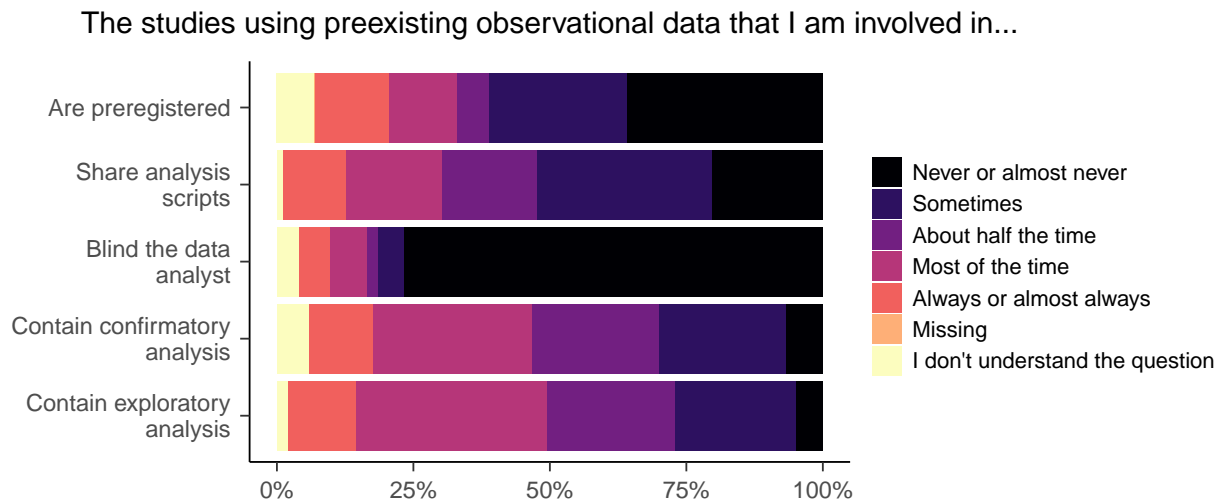
---

<sup>3</sup> The survey defined trustworthy as: “meaning that the results and conclusions of the publications are valid, reliable, rigorous, and accurate. That they merit trust.”

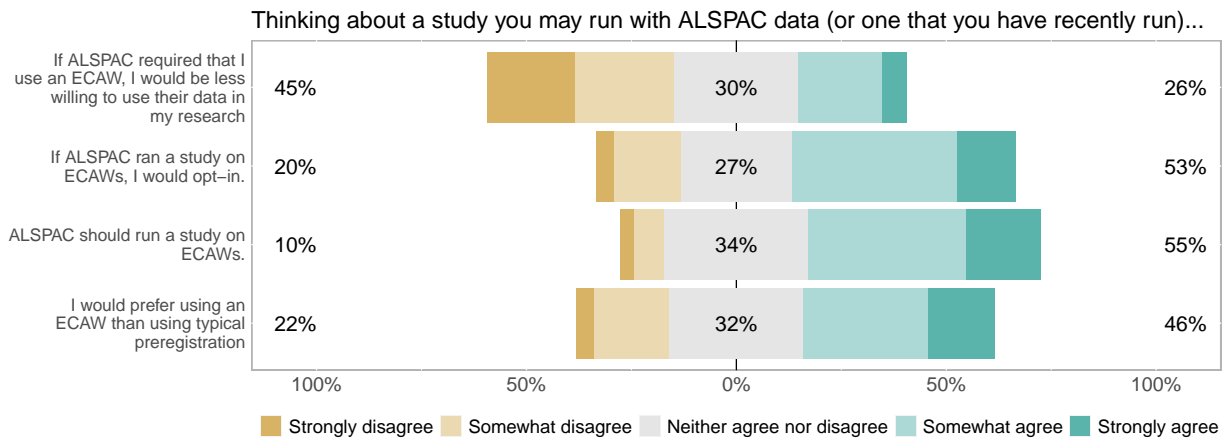
<sup>4</sup> The survey defined reproducible “in the sense that other researchers re-analysing the data with the same research question would produce similar results.”



*Figure 1. Responses to the survey questions on trustworthiness and reproducibility of observational research with preexisting data and ECAWs.* The survey defined trustworthy as “meaning that the results and conclusions of the publications are valid, reliable, rigorous, and accurate. That they merit trust”. The survey defined reproducible “in the sense that other researchers re-analysing the data with the same research question would produce similar results.” For each item, the number to the left of the data bar indicates the combined percentage for the responses depicted in any shade of brown/orange. The number in the center of the data bar (gray) indicates the percentage of neutral responses. The number to the right of the data bar indicates the combined percentage for the responses depicted in any shade of green. For the bottom panel we excluded the missing responses ( $n = ;$ ) and responses of “I don’t understand the question” ( $n = 3; 2$ ).



*Figure 2.* Responses to survey questions about the research practices of participants.



*Figure 3.* Responses to survey questions about using ECAWs. For the 4 questions we excluded missing values ( $n = 0; 0; 0; 0$ ), responses of “I don’t understand the question” ( $n = 0; 4; 1; 1$ ), and responses of “Unsure” ( $n = 2; 5; 6; 8$ ).

benefits and drawbacks of ECAWs, related research practices, and general comments. These questions received a total of (92) responses from (55) unique respondents. A complete list of responses are viewable in the open data [LINK]. We synthesized the response to open-ended questions into the 9 topics on the left side of this table. We divide these into three sections: (i) concerns about the acceptability of ECAWs, (ii) concerns that ECAWs will not have their intended impact, and (iii) alternative interventions that may achieve similar goals as typical preregistration and ECAWs. On the right side of the table, we provide a reflection on each topic.

## Exploratory analyses

## Discussion

## Acknowledgements



## References

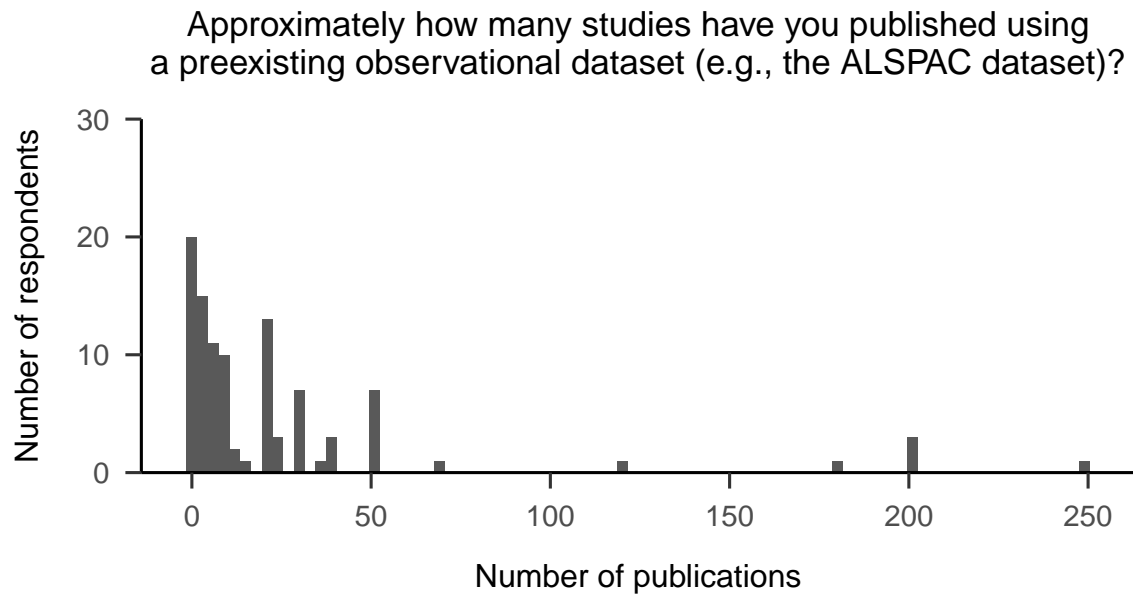


Figure 4. Caption goes here...

Table 1

*What programming language or software do you use for your analyses of preexisting observational data?*

<b>Programming language</b>	<b>N</b>	<b>Percentage of respondents</b>
R	65	63
Stata	48	47
SPSS	17	17
SAS	15	15
Python	6	6
Mplus	3	3
Bash	2	2
MATLAB	1	1
Nextflow	1	1
plink2	1	1

Compared to what you think of as a typical researcher who uses preexisting observational data in your field, how concerned are you with research trustworthiness, bias, rigour, and reproducibility ...

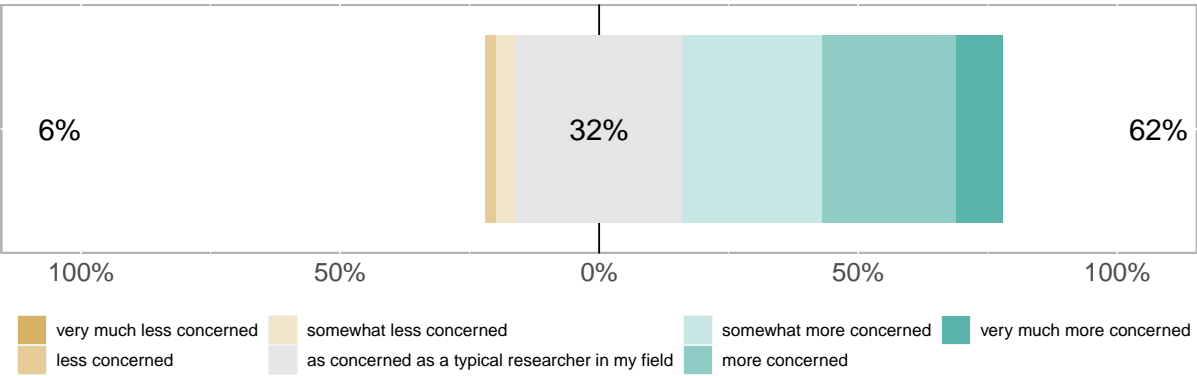


Figure 5. Caption goes here...