¹ Analytical choices for analyzing multidimensional behavior - many analyst test hypotheses

² about human speech.

³ First Author#, Second Author#, ...#, & Last Author#

⁴ $^{1}$ #

⁵ ... ...

⁶ Author Note

⁷ Add complete departmental affiliations for each author here. Each new line herein

⁸ must be indented, like this line.

⁹ Enter author note here.

¹⁴ Correspondence concerning this article should be addressed to First Author, Postal

¹⁵ address. E-mail: my@email.com

Abstract

One or two sentences providing a **basic introduction** to the field, comprehensible to a scientist in any discipline.

Two to three sentences of **more detailed background**, comprehensible to scientists in related disciplines.

One sentence clearly stating the **general problem** being addressed by this particular study.

One sentence summarizing the main result (with the words "**here we show**" or their equivalent).

Two or three sentences explaining what the **main result** reveals in direct comparison to what was thought to be the case previously, or how the main result adds to previous knowledge.

One or two sentences to put the results into a more **general context**.

Two or three sentences to provide a **broader perspective**, readily comprehensible to a scientist in any discipline.

*Keywords:* crowdsourcing science, data analysis, scientific transparency, speech, acoustic analysis

Word count: X

Analytical choices for analyzing multidimensional behavior - many analyst test hypotheses about human speech.

## Introduction

In order to effectively accumulate knowledge, science needs to (i) produce data that can be replicated using the original methods and (ii) arrive at robust conclusions substantiated by the data. In recent coordinated efforts to replicate published findings, the scientific disciplines have uncovered surprisingly low success rates for (i) (e.g., Open Science Collaboration, 2015; Camerer et al., 2018) leading to what is now referred to as the *replication crisis.* Beyond the difficulties of replicating scientific findings, a growing body of evidence suggests that the theoretical conclusions drawn from data are often variable even when researchers have access to reliable data (REFS). The latter situation has been referred to as the *inference crisis* (Rotello, Heit & Dubé 2015, Starns et al. 2019) and is, among other things, rooted in the inherent flexibility of data analysis (often referred to as researcher degrees of freedom: Simmons, Nelson, & Simonsohn, 2011, Gelman & Loken 2013). Data analysis involves many different steps, such as inspecting, organizing, transforming, and modeling the data, to name a few. Along the way, different methodological and analytical choices need to be made, all of which may influence the final interpretation of the data. These researcher degrees of freedom are both a blessing and a curse at the same time.

They are a blessing because they afford us the opportunity to look at nature from different angles, which, in turn, allows us to make important discoveries and generate new hypothesis (e.g. Box 1976, Tukey 1977, de Groot 2014). They are a curse because idiosyncratic choices can lead to categorically different interpretations, which eventually find their way into the publication record where they are taken for granted (Simmons et al. 2011). Recent projects have shown that the variability between different data analysts is vast. This variability can lead independent researchers to draw vastly different conclusions about the same dataset (e.g. Silberzahn et al. 2018, Starns et al. 2019, Botvinik-Nezer et al., 2020).

60  These projects, however, might still underestimate the extent to which analysts vary because

61  data analysis is not merely restricted to statistical inference of datasets. Human behavior is

62  complex and offers many ways to be translated into numbers. This is particularly true for

63  fields that draw conclusions about human behavior and cognition from multidimensional

64  data like speech or video data. In fields working on human speech production, for example,

65  researchers need to make numerous decisions about what to measure and how to measure it.

66  This is not trivial given the temporal extension of the acoustic signal and its complex

67  structural composition. Not only can decisions about operationalizing the raw data influence

68  downstream decisions about statistical modelling, but statistical results can also lead

69  researchers to go back and revise earlier decisions about the processing of the raw data.

70      In this article, we investigate the variability in analytic choices when many analyst

71  teams analyze the same speech production data. We explore the interaction between analytic

72  choices at the stage of the operationalization of raw data and subsequent statistical

73  modelling. Specifically, we report the impact of the analytic pipeline on research results

74  obtained by XX teams that gained access to the same set of acoustic recordings to answer

75  the same research question.

76  **Researcher degrees of freedom**

77      Data analysis comes with many decisions like how to operationalize a given

78  phenomenon or behavior, what data to submit to statistical modelling and which to exclude,

79  what models to use or what inferential decision procedure to apply. However, if these

80  decisions during data analysis are not specified in advance, we might stumble upon seemingly

81  meaningful patterns in the data that are merely statistical flukes. This can be problematic

82  because to err is human.

83      We have evolved to filter the world in irrational ways (e.g., Tversky and Kahneman

84  1974), seeing coherent patterns in randomness (Brugger 2001), convincing ourselves of the

85  validity of prior expectations ("I knew it", Nickerson 1998), and perceiving events as being

86  plausible in hindsight ("I knew it all along", Fischhoff 1975). In connection with an academic

87  incentive system that rewards certain discovery processes more than others (Sterling 1959,

88  Koole & Lakens 2012), we often find ourselves exploring many possible analytical pipelines,

89  but only reporting a select few. This issue is particularly amplified in fields in which the raw

90  data lend themselves to flexible measurement (Roettger 2019). Combined with a wide

91  variety of methodological and theoretical traditions as well as varying levels of statistical

92  training across fields and subfields, the inherent flexibility of data analysis might lead to a

93  vast plurality of analytic approaches.

94      Consequently, if methodologists are correct (e.g. Simmons et al. 2011, Gelman & Loken

95  2013), there are many published papers that present overconfident interpretations of their

96  data based on idiosyncratic analytic strategies. These interpretation are either associated

97  with an unknown amount of uncertainty or lend themselves to alternative interpretation if

98  analyzed differently. However, instead of being critically evaluated, scientific results often

99  remain unchallenged in the publication record. Despite recent efforts to improve

100  transparency and reproducibility (REFS) and freely available and accessible infrastructures

101  such as provided by the Open Science Framework (osf.io, ADD), critical reanalyses of

102  published analytic strategies are uncommon because data sharing is still rare (Wicherts,

103  Borsboom, Kats, & Molenaar, 2006, RECENT REF).

104      While this issue has been widely discussed both from a conceptual point of view

105  (Simmons et al. 2011, Wagenmakers et al. 2012, Nosek and Lakens 2014) and its application

106  in individual scientific fields (e.g. Wichert et al. 2015, Charles et al. 2019, Roettger 2019),

107  there are still many unknowns regarding the extent of analytical plurality in practice. Recent

108  collaborative attempts have started to shed light on how different analysts tackle the same

109  data set and have revealed, not surprisingly, a large amount of variability.

**Crowdsourcing alternative analyses**

In a collaborative effort, Silberzahn et al. (2018) let twenty-nine independent analysis teams address the same research hypothesis Analytic approaches and consequently the results varied widely between teams. Sixty-nine percent of the teams found support for the hypothesis, and 31% did not. Out of the 29 analytical strategies, there were 21 unique combinations of covariates. Importantly, the observed variability was neither predicted by the team's preconceptions about the phenomenon under investigation nor by peer ratings of the quality of their analyses. The authors results suggest that analytic plurality is a fact of life and not driven by different levels of expertise or bias. Similar crowd-sourced studies recruiting independent analyst teams showed similar results.

SUM UP: Neuroscience

Cognitive Modelling

Clinical

Predictive models

While these projects show a large degree of analytical flexibility with impactful consequences, they dealt with flexibility in inferential or computational modelling. In these studies the datasets were fixed and data collection or extraction could not be changed.

However, in many fields the primary raw data is a complex signal that needs to be operationalized according to the research question. In social sciences, the raw observations correspond to recorded human behavior. In many cases, the behavior is recorded and stored as a complex visual and/or acoustic signal that is temporally extended and exhibits a complex structure. Decisions about how to operationalize a theoretically relevant aspect of that behavior or the underlying cognitive processes might interact with downstream decisions about statistical modelling and vice verse. To understand how analytical flexibility manifests itself in a scenario where a complex decisions procedure is involved in operationalizing and measuring complex signals, the present paper looks at an

136  experimentally elicited speech data set

## Operationalizing speech

138          RELEVANCE OF SPEECH PRODUCTION RESEARCH FOR COGSCI, AI, etc.

139          In order to understand speech, listeners have to map a continuous, transient signal

140  onto discrete meanings. Take for example the sentence: "The president refused to concede".

141  This sentence stretches over hundreds of milliseconds and contains several layers of

142  linguistically relevant units such as words, syllables and individual sounds. Thus, it offers a

143  considerable number of perspectives and decisions along the data analysis pipeline.

144          IMAGE SHOWING A WAVE FORM WITH DIFFERENT DOMAINS (temporal)

145  AND STRUCTURAL CUES (e.g. f0, dur, int) MAPPING ONTO FUNCTIONAL

146  CONTRASTS

147          When conducting a study on speech production, the first important analytic decision to

148  test a hypothesis relates to operationalization of the relevant behavior, i.e. how to measure

149  the phenomenon of interest. For example, how do we measure whether two sounds are

150  acoustically identical or not (e.g. "bear" vs. "pear"), whether one word is more prominent

151  than others ("The president REFUSES to concede"), or whether two discourse functions are

152  produced differently ("The president refuses to concede." vs. "The president refuses to

153  concede?")? In other words, how do we quantitatively capture relevant features of speech?

154          This is not trivial. Speech categories are inherently multidimensional and vary through

155  time. The acoustic parameters for one category are usually asynchronous, i.e. appear at

156  different points of time in the unfolding signal and overlap with parameters for other

157  categories (e.g. Jongman et al., 2000; Lisker, 1986; Summerfield, 1984; Winter, 2014). For

158  example, the distinction between voiced and voiceless stops in English (i.e. /b/ and /p/ in

159  "bear" vs. "pear") can be manifested by many different acoustic measures (Lisker, 1977).

160  Even temporally dislocated acoustic parameters correlate with this lexical contrast. For

161  example, in the English words "led" versus "let", voicing correlates can be found in the

162  acoustic manifestation of the initial /l/ of the word (Hawkins & Nguyen, 2004).

163      The apparent multiplicity of phonetic cues grows exponentially if we look at larger

164  temporal windows as is the case for suprasegmental aspects of speech. Studies investigating

165  acoustic correlates of word stress (e.g. the difference between "insight" and "incite"), for

166  example, have been using many different measurements, including temporal characteristics

167  (duration of certain segments or subphonemic intervals), spectral characteristics (intensity

168  measures, formants, and spectral tilt), and measurements related to fundamental frequency

169  (f0) (e.g. Gordon & Roettger, 2017).

170      Looking at even larger domains, the expression of pragmatic functions can be

171  expressed by a variety of structurally different acoustic cues which can be distributed

172  throughout the entire utterance. Discourse functions are systematically expressed by

173  multiple local pitch modulations differing in their position, shape, and alignment

174  (e.g. Niebuhr et al., 2011). They can also be expressed by global or local pitch modulations,

175  as well as acoustic information within the temporal or spectral domain (e.g. van Heuven &

176  van Zanten 2005). All of these phonetic parameters are potential manifestations of

177  underlying functional contrasts, like speaker's intentions, levels of arousal or social identity.

178      When testing hypotheses on speech production data, researchers need to make many

179  decisions. The larger the functional domain (e.g. are we interested in lexical items or in

180  whole utterances), the higher the number of possible operationalizations. These decisions are

181  usually made prior to any statistical analysis, but can possibly be revised a posteriori in light

182  of downstream analytic decisions. To probe the variability in data analysis pipelines across

183  independent analytic teams, we provided researchers with an experimentally elicited speech

184  corpus that looked at a functional contrast that is potentially manifested across the whole

185  utterance.

**The data set - The prosody of redundant modifiers**

Our data set was collected in order to answer the following research question: Do speakers acoustically modify utterances to signal atypical word combinations? (e.g. "a blue banana" vs. "a yellow banana")?

Referring is one of the most basic and prevalent uses of language and one of the most widely researched areas in language science. It is an open question how speakers choose a referential expression when they want to refer to a specific entity like a banana. The context within which an entity occurs (i.e., with other non-fruits, other fruits, or other bananas) plays a large part in determining the choice of referential expression. Generally, speakers aim to be as informative as possible to uniquely establish reference to the intended object, but they are also resource efficient in that they avoid redundancy (Grice 1975). Thus one would expect the use of a modifier, for example, only if it is necessary for disambiguation. For instance, one might use the adjective "yellow" to describe a banana in a situation in which there is a yellow and a less ripe green banana available, but not when there is only one banana to begin with.

Despite this coherent idea that speakers are both rational and efficient, there is much evidence that speakers are often overinformative: Speakers use referring expressions that are more specific than strictly necessary for the unambiguous identification of the intended referent (Sedivy 2003, Westerbeek et al. 2015, Rubio-Fernandez 2016), which has been argued to facilitate object identification and making communication between speakers and listeners more efficient (Arts et al. 2011, Paraboni et al. 2007, Rubio-Fernandez 2016). Recent findings suggest that the utility of a referring expression depends on how good it is for a listener (compared to other referring expressions) to identify a target object. For example, Degen et al. (2020) showed that modifiers that are less typical for a given referent (e.g. a blue banana) are more likely to be used in an overinformative scenario (e.g. when there is just one banana).

²¹²      This account, however, has mainly focused on content selection (Gatt et al. 2013),

²¹³ i.e. whether a certain referential expression is chosen or not, ignoring the fact that speech

²¹⁴ communication is much richer. Even looking at morphosyntactically identical expressions,

²¹⁵ speakers can modulate these expressions via suprasegmental acoustic properties like

²¹⁶ temporal and spectral modifications of the segments involved (e.g. Ladd 2008). Most

²¹⁷ prominently, languages use intonation to signal discourse relationships between referents

²¹⁸ (among other functions). Intonation marks discourse-relevant referents for being new or

²¹⁹ given information to guide listeners' interpretation of incoming messages. In many languages,

²²⁰ speakers can use particular pitch movements to signal whether a referent has already been

²²¹ mentioned and is therefore referred back to, or a referent is newly introduced into the

²²² discourse. Many languages use intonation in order to signal if a referent is contrasting with

²²³ one or more alternatives that are relevant to the current discourse. Content selection aside,

²²⁴ in a scenario in which a speaker wants to refer to a banana when there is also a pear on the

²²⁵ table, the speaker would most likely produce a high rising pitch accent on "banana" to

²²⁶ indicate the contrastive nature of the *noun.* In a scenario in which the speaker wants to refer

²²⁷ to a yellow banana when there is also a less ripe green banana on the table, the speaker

²²⁸ would most likely produce a high rising pitch accent on "yellow" to indicate the contrastive

²²⁹ nature of the *modifier.* In addition to a pitch accent, elements that are new and/or

²³⁰ contrastive are often produced with additional suprasegmental prominence, i.e. segments are

²³¹ hyperarticulated, resulting in longer, louder and more clearly articulated acoustic targets.

²³²      INFORMATION ABOUT THE DATA SET AND EXP DESIGN

²³³ **Research questions**

²³⁴                               **Methods (mostly copy-paste from Evo-RR)**

²³⁵      We are closely following the methodology proposed by Parker et al. (Stage 1

²³⁶ in-principle accepted).

237     This project involves a series of steps (X-X) that begins with recruiting independent

238  groups of scientists to analyze the data, continuing through allowing the scientists to analyze

239  the data as they see fit, generating peer review ratings of the analyses (based on methods,

240  not results), evaluating the variation among the different analyses, and producing the final

241  manuscript. We estimate that this process, from the time of an in-principle acceptance of

242  this Stage 1 Registered Report, will take XX months (Table X). The factor most likely to

243  delay our time line is the rate of completion of the original set of analyses by independent

244  groups of scientists.

**Step 1: Recruitment and Initial Survey of Analysts**

246     Initiating authors (SC, JC, TR) created a publicly available document providing a

247  general description of the project (LINK) and a short prerecorded slide show that

248  summarizes the study and research question in order to increase accesibilty to potential

249  analysts (LINK). The project will be advertised via Social Media, using mailing lists for

250  linguistic and psychological societies (full scope of these lists is not fixed but will include

251  LIST OF LISTS), and via word of mouth. The target population is active speech science

252  researchers with a graduate degree (or currently studying for a graduate degree) in a relevant

253  discipline. Researchers can choose to work independently or in a small team. For the sake of

254  simplicity, we refer to single researcher or small teams as "analysis teams".

255     Recruitment for this project is ongoing but we aim for a minimum of XX analysis

256  teams independently evaluating each dataset (see sample size justification below). We will

257  simultaneously recruit volunteers to peer-review the analyses conducted by the other

258  volunteers through the same channels. Our goal is to recruit a similar number of

259  peer-reviewers and analysts, and to ask each peer reviewer to review a minimum of four

260  analyses. If we are unable to recruit at least half the number of reviewers as analysis teams,

261  we will ask analysts to serve also as reviewers (after they have completed their analyses). All

262  analysts and reviewers will share co-authorship on this manuscript and will participate in the

263 collaborative process of producing the final manuscript. All analysts will sign a consent

264 (ethics) document (LINK).

265        We identified our minimum number of analyses per data set by considering the number

266 of effects needed in a meta-analysis to generate an estimate of heterogeneity $(\tau^2)$ with a 95%

267 confidence interval that does not encompass zero. This minimum sample size is invariant

268 regardless of $\tau^2$. This is because the same t-statistic value will be obtained by the same

269 sample size regardless of variance $(\tau^2)$. We see this by first examining the formula for the

270 standard error, SE for variance, $(\tau^2)$ or $SE(\tau^2)$ assuming normality in an underlying

271 distribution of effect sizes (Knight 2000):

$$SE(\tau^2) = \sqrt{\frac{2\tau^4}{(n-1)}}$$

272        and then rearranging the above formula to show how the t-statistic is independent of

273 $\tau^2$, as seen below.

$$t = \frac{\tau^2}{SE(\tau^2)} = \sqrt{\frac{(n-1)}{2}}$$

274        We then find a minimum n = 12 according to this formula.

**Step 2: Primary Data Analyses**

276        Analysis teams will register and answer a demographic and expertise survey (LINK).

277 The survey collects information on the analysts current position and self-estimated breadth

278 and level of statistical expertise. We will then provide teams with the acoustic data set and

279 request that they answer the following research question:

280        Do speakers acoustically modify utterances to signal atypical word combinations?

281        Once their analysis is complete, they will answer a structured survey (LINK),

282 providing analysis technique, explanations of their analytical choices, quantitative results,

283 and a statement describing their conclusions. They will also upload their analysis files

284 (including the additionally derived data and text files that were used to extract and

285 preprocess the acoustic data), their analysis code (if applicable), and a detailed journal-ready

286 statistical methods section.

**Step 3: Peer Reviews of Analyses**

288        At a minimum, each analysis will be evaluated by four different reviewers, and each

289 volunteer peer-reviewer will be randomly assigned to methods sections from at least four

290 analyst teams (the exact number will depend on the number of analysis teams and peer

291 reviewers recruited). Each peer reviewer will register and answer a demographic and

292 expertise survey identical to that asked of the analysts. Reviewers will evaluate the methods

293 of each of their assigned analyses one at a time in a sequence determined by the initiating

294 authors (SC, JC, TR). The sequences will be systematically assigned so that, if possible,

295 each analysis is allocated to each position in the sequence for at least one reviewer. For

296 instance, if each reviewer is assigned four analyses to review, then each analysis will be the

297 first analysis assigned to at least one reviewer, the second analysis assigned to another

298 reviewer, the third analysis assigned to yet another reviewer, and the fourth analysis

299 assigned to a fourth reviewer. Balancing the order in which reviewers see the analyses

300 controls for order effects, e.g. a reviewer might be less critical of the first methods section

301 they read than the last. The process for a single reviewer will be as follows. First, the

302 reviewer will receive a description of the methods of a single analysis. This will include the

303 narrative methods section, the analysis team's answers to our survey questions regarding

304 their methods, including analysis code, and the data set. The reviewer will then be asked, in

305 an online survey (LINK), to rate both the acoustic analysis and the statistical analysis on a

306 scale of 0-100 based on these prompts:

307        "Rate the overall appropriateness of the acoustic analysis to answer the research

308 question with the available data. To help you calibrate your rating, please consider the

following guidelines:

- 100. A perfect analysis with no conceivable improvements from the reviewer

- 75. An imperfect analysis but the needed changes are unlikely to dramatically alter final interpretation

- 50. A flawed analysis likely to produce either an unreliable estimate of the relationship or an over-precise estimate of uncertainty

- 25. A flawed analysis likely to produce an unreliable estimate of the relationship and an over-precise estimate of uncertainty

- 0. A dangerously misleading analysis, certain to produce both an estimate that is wrong and a substantially over-precise estimate of uncertainty that places undue confidence in the incorrect estimate.

*Please note that these values are meant to calibrate your ratings. We welcome ratings of any number between 0 and 100."

After providing this rating, the reviewer will then be provided with a series of text boxes and the following prompts:

"Please explain your ratings of this analysis.

Please evaluate the selection of acoustic features.

Please evaluate the measurement of acoustic features.

Please evaluate the choice of statistical analysis type.

Please evaluate the process of choosing variables for and structuring the statistical model.

Please evaluate the suitability of the variables included in (or excluded from) the statistical model.

Please evaluate the suitability of the structure of the statistical model.

Please evaluate choices to exclude or not exclude subsets of the data.

Please evaluate any choices to transform data (or, if there were no transformations, but you think there should have been, please discuss that choice)."

After submitting this review, a methods section from a second analysis will then be made available to the reviewer. This same sequence will be followed until all analyses allocated to a given reviewer have been provided and reviewed. After providing the final review, the reviewer will be simultaneously provided with all four (or more) methods sections that reviewer has just completed reviewing, the option to revise their original ratings, and a text box to provide an explanation. The invitation to revise the original ratings will be as follows: "If, now that you have seen all the analyses you are reviewing, you wish to revise your ratings of any of these analyses, you may do so now." The text box will be prefaced with this prompt: "Please explain your choice to revise (or not to revise) your ratings."

**Step 4: Evaluate Variation**

Initiating authors (SC, JC, TR) will conduct the analyses outlined in this section. We will describe the variation in model specification in several ways: First, we will calculate summary statistics describing variation among analysis, including the nature and number of acoustic measures, the operationalization and the temporal domain of measurement, the nature and number of model parameters for both fixed and random effects [if applicable], as well as the mean, standard deviation and range of effect sizes reported.

We will summarize the variability in standardized effect sizes and predicted values of dependent variables among the individual analyses using standard random effects meta-analytic techniques. We anticipate that the majority of statistical analyses will be expressible as a (generalized) linear regression model.

ADD FORMULA

First, we will derive standardized effect sizes from each individual analysis. Since we anticipate that researchers use multi-level linear regression models, common effect size

358    measures such as Cohen's d are inappropriate. In multi-level models the variance

359    components are partitioned into multiple sources of variation (e.g. varying intercept of

360    speakers and items, varying by-predictor slopes, etc.). We thus will take into account all of

361    the variance sources of the models (Hedges, 2007) and derive the index $\delta_t$ (where t stands for

362    "total" variance), which is calculated by the estimated difference between group means ($\beta$),

363    divided by the square root of the sum of all variance components as formulated below:

$$\delta_t = \frac{\beta}{\sqrt{\sum_{i=1}^{n} \sigma_i^2}}$$

364    Where $_i$ refers to the individual variance components.

365    Variation in the resulting effect sizes can be due to many different sources, one of which

366    is the model architecture and the specification of random effect structure. Hypothetically,

367    two analyst teams can arrive at two different effect sizes, even if they have made the exact

368    same measurement, processed the data in the exact same way, and used the same predictor

369    combination, but, crucially, differ in their random effect structure (e.g. one team assumes

370    random intercepts only, the other team uses random slopes). Thus, additionally, we will run

371    all analyses with a prespecified maximal random effect structure for all predictors.

372    Upon calculating a standardized effect size and standard error for each analysis, the

373    initiating authors will then fit a cross-classified Bayesian meta-analysis on the analyst team

374    data using the multilevel regression model described below:

$$\delta_t \sim \text{Normal}(\theta_i, \sigma_i = \text{se}_i)$$

$$\theta_i \sim \text{Normal}(\mu, \tau)$$

$$\mu \sim \text{Normal}(0, 1)$$

$$\tau \sim \text{HalfCauchy}(0, 1)$$

375    Effect size ($\delta_t$) is the outcome variable and the number of post-hoc changes and the

376    number of models fit will be included as population-level effects (i.e., fixed effects). The

377 likelihood of the outcome variable is assumed to be normally distributed. Analysis teams and

378 reviewers will be included as group-level effects (i.e., random effects). For all population-level

379 parameters, the model will include regularizing, weakly informative priors (Gelman, 2017),

380 which are normally distributed and centered at 0 with a standard deviation of 1. A cauchy

381 prior set at 0 with scale 1 will be used for $\tau$. We will fit the model with 4000 iterations (2000

382 warm-up) and Hamiltonian Monte-Carlo sampling of the posterior distribution is carried out

383 using 4 chains distributed across 4 processing cores. The analysis will be conducted in R (R

384 core team, 2020) and fit using `stan` (Stan, 2019) via the R package `brms` (Bürkner, 2019).

385     The pooled estimate of the meta-analytic model will be used in a series of descriptive

386 analyses that serve the purpose of describing how individual teams vary from each other.

387 Specifically, we will explain how team effect sizes deviate from the pooled estimate based on

388 a measure of the uniqueness of the set of variables included in each analysis. We will use

389 Sorensen's Similarity Index (SSI) to derive a "uniqueness" score. The SSI is an index

390 typically used in ecology research to compare species composition across sites. For our

391 purposes, we will treat variables as species and individual analyses as sites. In order to

392 generate an SSI for each analysis team, we will calculate the average of all pairwise

393 Sorensen's values for all pairs of analyses using the betapart package (Baselga et al. 2018) in

394 R. We achieve this using the following formula:

$$\beta_{Sorensen} = \frac{(b+c)}{(2a+b+c)}$$

395     where, given a pair of models, a is the number of variables common to both, b is the

396 number of variables that occur in the first model but not in the second, and c is the number

397 of variables that occur in the second model but not in the first. We then will use the

398 per-model average Sorensen's index value as an independent variable to predict the deviation

399 score in a general linear model, with no random effect since each analysis is included only

400 once, in the stats package in R (R_Core_Team 2019):

WRITE ABOUT HOW TO HANDLE GAMS

We will publicly archive all relevant data, code, and materials on the Open Science Framework (ADD LINK). Archived data will include the original data sets distributed to all analysts, any edited versions of the data analyzed by individual groups, and the data we analyze with our meta-analyses, which include the effect sizes derive from separate analyses, the statistics describing variation in model structure among analyst groups, and the anonymized answers to our surveys of analysts and peer reviewers. Similarly, we will archive both the analysis code used for each individual analysis and the code from our meta-analyses. We will also archive copies of our survey instruments from analysts and peer reviewers.

Our rules for excluding data from our study are as follows. We will exclude from our synthesis any individual analysis submitted after we have completed peer review or those unaccompanied by analysis files that allow us to understand what the analysts did. We will also exclude any individual analysis that does not produce an outcome that can be interpreted as an answer to our primary question.

ADD GAM EXCLUSION

DESCRIBE META-NALYSIS

We will assess the extent to which deviations from the meta-analytic mean by individual effect sizes can be explained by the following predictors:

The peer ratings (1-100) of each analysis, a measurement of the "uniqueness" of the set of model parameters included in each analysis (see below), whether or not analysis teams have changed their analyses in a post hoc fashion (EXPLAIN ABOVE, binary), the number of analyses performed before the final model was performed (as self-estimated by the teams).

The uniqueness of model parameters will be assessed using the Sorensen's Similarity Index (an index typically used to compare species composition across sites). Here we treat model parameters as species and individual analyses as sites. To generate an individual Sorensen's value for each analysis requires calculating the pairwise Sorensen's value for all

pairs of analyses (of the same data set), and then taking the average across these Sorensen's

values for each analysis. We will calculate the Sorensen's index values using the betapart

package (Baselga et al. 2018) in R:

$$\beta_{Sorensen} = \frac{b + c}{2_a + b + c}$$

where $a$ is the number of variables common to both models, $b$ is the number of

variables that occur in the first model but not in the second and $c$ is the number of variables

that occur in the second model but not in the first.

The individual posterior mean for each analysis of each team serves as the dependent

variable in these analyses.

We wish to quantify relationships among variables, but we have no a priori expectation

of effect size and we will not make dichotomous decisions (such as statistical significance).

**Step 6: Collaborative Write-Up of Manuscript**

Analysts and initiating authors will discuss the limitations, results, and implications of

the study and collaborate on writing the final manuscript for review as a stage-2 Registered

Report.

## References

Camerer, C. F., Dreber, A., Holzmeister, F., Ho, T.-H., Huber, J., Johannesson, M., ...
        others. (2018). Evaluating the replicability of social science experiments in nature
        and science between 2010 and 2015. *Nature Human Behaviour*, *2*(9), 637–644.
        https://doi.org/10.1038/s41562-018-0399-z

Open Science Collaboration. (2015). Estimating the reproducibility of psychological science.
        *Science*, *349*(6251). https://doi.org/10.1126/science.aac4716