

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

First Author[#], Second Author[#], ...[#], & Last Author[#]

¹ #

... ...

Author Note

Add complete departmental affiliations for each author here. Each new line herein
must be indented, like this line.

Enter author note here.

The authors made the following contributions. First Author: Conceptualization,
Writing - Original Draft Preparation, Writing - Review & Editing; Second Author: Writing -
Review & Editing; ...: Writing - Review & Editing; Last Author: Writing - Review &
Editing.

Correspondence concerning this article should be addressed to First Author, Postal
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 **broad 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) (Camerer et al., 2018; e.g., Open Science Collaboration, 2015) 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).

These projects, however, might still underestimate the extent to which analysts vary because data analysis is not merely restricted to statistical inference of datasets. Human behavior is complex and offers many ways to be translated into numbers. This is particularly true for fields that draw conclusions about human behavior and cognition from multidimensional data like audio or video data. In fields working on human speech production, for example, researchers need to make numerous decisions about what to measure and how to measure it. This is not trivial given the temporal extension of the acoustic signal and its complex structural composition. Not only can decisions about measuring the signal influence downstream decisions about statistical modelling, but statistical results or modelling issues can also lead researchers to go back and revise earlier decisions about the measuring process itself.

In this article, we investigate the variability in analytic choices when many analyst teams analyze the same speech production data, a process that involves both decisions regarding the operationalization of a complex observed signal and decisions regarding the statistical modelling. Specifically, we report the impact of the analytic pipeline on research results obtained by XX teams who gained access to the same set of acoustic recordings in order to answer the same research question.

Researcher degrees of freedom

Data analysis comes with many decisions like how to measure a given phenomenon or behavior, what data to submit to statistical modelling and which to exclude in the final analysis, what models to use or what inferential decision procedure to apply. However, if these decisions during data analysis are not specified in advance, we might stumble upon seemingly meaningful patterns in the data that are merely statistical flukes. This can be problematic because to err is human.

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

1974), seeing coherent patterns in randomness (Brugger 2001), convincing ourselves of the validity of prior expectations (“I knew it,” Nickerson 1998), and perceiving events as being plausible in hindsight (“I knew it all along,” Fischhoff 1975). In connection with an academic incentive system that rewards certain discovery processes more than others (Sterling 1959, Koole & Lakens 2012), we often find ourselves exploring many possible analytical pipelines, but only reporting a select few. This issue is particularly amplified in fields in which the raw data lend themselves to many possible ways to measure (Roettger 2019). Combined with a wide variety of methodological and theoretical traditions as well as varying levels of statistical training across subfields, the inherent flexibility of data analysis might lead to a vast plurality of analytic approaches that can lead to different scientific conclusions. Consequently, there might be many published papers that present overconfident interpretations of their data based on idiosyncratic analytic strategies (e.g. Simmons et al. 2011, Gelman & Loken 2013). These interpretations are either associated with an unknown amount of uncertainty or lend themselves to alternative interpretation if analyzed differently. However, instead of being critically evaluated, scientific results often remain unchallenged in the publication record. Despite recent efforts to improve transparency and reproducibility (REFS) and freely available and accessible infrastructures such as provided by the Open Science Framework (osf.io, ADD), critical reanalyses of published analytic strategies are still not very common because data sharing remains rare (Wicherts, Borsboom, Kats, & Molenaar, 2006, RECENT REF).

While this issue has been widely discussed both from a conceptual point of view (Simmons et al. 2011, Wagenmakers et al. 2012, Nosek and Lakens 2014) and its application in individual scientific fields (e.g. Wichert et al. 2015, Charles et al. 2019, Roettger 2019), there are still many unknowns regarding the extent of analytical plurality in practice. Recent collaborative attempts have started to shed light on how different analysts tackle the same data set and have revealed 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 measurement could not be changed.

However, in many fields the primary raw data are complex signals that need to be operationalized according to the research question. In social sciences, the raw observations correspond to human behavior, sometimes measured as a complex visual or acoustic signal. Decisions about how to measure a theoretically construct related to that behavior or the underlying cognitive processes might interact with downstream decisions about statistical modelling and vice versa.

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 experimentally elicited speech data set

Operationalizing speech

One of the earliest analytical decisions a researcher has to make when conducting a study on speech production is choosing how to measure the phenomenon of interest, i.e. how to operationalize it. Take for example the following sentence: “Sheila says Pat is clever.” Now, imagine we are interested in measuring the difference between that sentence and the sentence “Sheila says *Mat* is clever.” How can we obtain a measure of dissimilarity between the first sound in *Pat* and the first sound in *Mat*? If we want to compare “Sheila says PAT is clever” to “Sheila says Pat is CLEVER,” how do we quantify the difference between the two utterances? What if we are interested in how the argumentative nature of “Sheila says Pat is clever” and the expression of surprise in “Sheila says Pat is clever?!” are conveyed in speech? How do we choose what to measure and how to measure it so that we can answer our research question and test our hypotheses? Given the continuous and transient nature of speech, identifying which speech features should be selected within which domain becomes a non-trivial task. Utterances stretch over hundreds of milliseconds and contain several levels of linguistically relevant units such as phrases, words, syllables, and individual sounds. The researcher is thus confronted with a considerable number of features and combinations thereof to choose from.

Speech categories are inherently multidimensional and dynamic: they consist of a cluster of features that are modulated over time. The acoustic signatures of one category are usually asynchronous, i.e. they appear at different time points in the unfolding signal, and overlap with the signatures of other categories (e.g. Jongman et al., 2000; Lisker, 1986; Summerfield, 1984; Winter, 2014). A classical example is the distinction between voiced and voiceless stops in English (i.e. /b/ and /p/ in *bear* vs *pear*). This voiced/voiceless contrast is manifested by many acoustic features which can differ depending on several factors, such as position of the consonant in the word and surrounding sounds (Lisker, 1977). Furthermore, correlates of the contrast can even be found away from the consonant, in temporally distant

speech units. The initial /l/ of the English words *led* and *let* is affected by the voicing of the final consonant (/t, d/) (Hawkins & Nguyen, 2004). The multiplicity of phonetic cues grows exponentially if we look at larger temporal windows as is the case for suprasegmental aspects of speech. Studies investigating acoustic correlates of word stress (e.g. the difference between *insight* and *incite*) have been using a wide variety of measurements, including temporal characteristics (duration of certain segments or sub-segmental intervals), spectral characteristics (intensity, formants, and spectral tilt), and measurements related to fundamental frequency (f0) (e.g. Gordon & Roettger, 2017).

Moving onto the expression of higher-level functions like information structure and discourse pragmatics, the relevant acoustic cues can be distributed throughout even larger domains, such as phrases and whole sentences. Differences in position, shape, and alignment of pitch modulations over multiple locations within a sentence are correlated with differences in discourse functions (e.g. Niebuhr et al., 2011). The latter can also be expressed by global vs local pitch modulations, as well as acoustic information within the temporal or spectral domain (e.g. van Heuven & van Zanten 2005). Extra- and para-linguistic information, like speaker’s intentions, levels of arousal or social identity, are also conveyed by broad-domain features, such as voice quality, speech rate, and rhythm.

When testing hypotheses on speech production data, researchers are faced with many choices and possibilities. The larger the functional domain (e.g. segments vs words vs utterances), the higher the number of conceivable operationalizations. Moreover, even the analysis of a single measure can be approached via an ever-increasing range of different statistical models, which further multiply the combinations of possible analytical choices. These decisions are usually made prior to any statistical analysis, but are at times revised a posteriori (i.e. after data collection and/or preliminary analyses) in light of unforeseen or surprising outcomes. To probe the variability in data analysis pipelines across independent researchers, we provided analytical teams with an experimentally elicited speech corpus and asked them to investigate acoustic differences related to a functional contrast that might be

187 manifested across the whole utterance.

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

189 Our data set was collected in order to answer the following research question: Do
190 speakers acoustically modify utterances to signal atypical word combinations? (e.g. “a blue
191 banana” vs. “a yellow banana”)? We are interested in the acoustic profile of referring
192 expression. Referring is one of the most basic and prevalent uses of language and one of the
193 most widely researched areas in language science. It is an open question how speakers choose
194 a referential expression when they want to refer to a specific entity like a banana. The
195 context within which an entity occurs (i.e., with other non-fruits, other fruits, or other
196 bananas) plays a large part in determining the choice of referential expression. Generally,
197 speakers aim to be as informative as possible to uniquely establish reference to the intended
198 object, but they are also resource efficient in that they avoid redundancy (Grice 1975). Thus
199 one would expect the use of a modifier, for example, only if it is necessary for
200 disambiguation. For instance, one might use the adjective “yellow” to describe a banana in a
201 situation in which there is a yellow and a less ripe green banana available, but not when
202 there is only one banana to begin with. Despite this coherent idea that speakers are both
203 rational and efficient, there is much evidence that speakers are often over-informative:
204 Speakers use referring expressions that are more specific than strictly necessary for the
205 unambiguous identification of the intended referent (Sedivy 2003, Westerbeek et al. 2015,
206 Rubio-Fernandez 2016), which has been argued to facilitate object identification and making
207 communication between speakers and listeners more efficient (Arts et al. 2011, Paraboni et
208 al. 2007, Rubio-Fernandez 2016). Recent findings suggest that the utility of a referring
209 expression depends on how good it is for a listener (compared to other referring expressions)
210 to identify a target object. For example, Degen et al. (2020) showed that modifiers that are
211 less typical for a given referent (e.g. a blue banana) are more likely to be used in an
212 over-informative scenario (e.g. when there is just one banana). This account, however, has

mainly focused on content selection (Gatt et al. 2013), i.e. whether a certain referential expression is chosen or not, ignoring the fact that speech communication is much richer. Even looking at morphosyntactically identical expressions, speakers can modulate these expressions via suprasegmental acoustic properties like temporal and spectral modifications of the segments involved (e.g. Ladd 2008). Most prominently, languages use intonation to signal discourse relationships between referents (among other functions). Intonation marks discourse-relevant referents for being new or given information to guide listeners' interpretation of incoming messages. In many languages, speakers can use particular pitch movements to signal whether a referent has already been mentioned and is therefore referred back to, or a referent is newly introduced into the discourse. Many languages use intonation in order to signal if a referent is contrasting with one or more alternatives that are relevant to the current discourse. Content selection aside, in a scenario in which a speaker wants to refer to a banana when there is also a pear on the table, the speaker would most likely produce a high rising pitch accent on 'banana' to indicate the contrastive nature of the *noun*. In a scenario in which the speaker wants to refer to a yellow banana when there is also a less ripe green banana on the table, the speaker would most likely produce a high rising pitch accent on 'yellow' to indicate the contrastive nature of the *modifier*. In addition to a pitch accent, elements that are new and/or contrastive are often produced with additional suprasegmental prominence, i.e. segments are hyperarticulated, resulting in longer, louder and more clearly articulated acoustic targets.

INFORMATION ABOUT THE DATA SET AND EXP DESIGN

Research questions

The present project examines the extent to which subjective choices by different researchers analyzing a complex speech data set affect the reported results. We are further interested in which factors affect researchers' final results.

Methods (mostly copy-paste from Evo-RR)

We are closely following the methodology proposed by Parker et al. (Stage 1 in-principle accepted) in terms of data collection. The analysis will substantially diverge from their approach (see §#. #)

This project involves a series of steps (X-X): First, we recruit independent groups of researchers to analyze the data. Second, We give researchers access to the speech corpus and let them analyze the data as they see fit. Third, we generate peer review ratings of the analyses (based on methods, not results). Forth, we evaluate the variation among the different analyses. And finally, we collaboratively produce the final manuscript. We estimate that this process, from the time of an in-principle acceptance of this Stage 1 Registered Report, will take XX months (Table X). The factor most likely to delay our time line is the rate of completion of the original set of analyses by independent groups of scientists.

Step 1: Recruitment and Initial Survey of Analysts

Initiating authors (SC, JC, TR) created a publicly available document providing a general description of the project (LINK) and a short prerecorded slide show that summarizes the study and research question in order to increase accesibilty to potential analysts (LINK). The project will be advertised via Social Media, using mailing lists for linguistic and psychological societies (full scope of these lists is not fixed but will include LIST OF LISTS), and via word of mouth. The target population is active speech science researchers with a graduate degree (or currently studying for a graduate degree) in a relevant discipline. Researchers can choose to work independently or in a small team. For the sake of simplicity, we refer to single researcher or small teams as ‘analysis teams.’

Recruitment for this project is ongoing but we aim for a minimum of XX analysis teams independently evaluating each dataset (see sample size justification below). We will simultaneously recruit volunteers to peer-review the analyses conducted by the other

volunteers through the same channels. Our goal is to recruit a similar number of peer-reviewers and analysts, and to ask each peer reviewer to review a minimum of four analyses. If we are unable to recruit at least half the number of reviewers as analysis teams, we will ask analysts to serve also as reviewers (after they have completed their analyses). All analysts and reviewers will share co-authorship on this manuscript and will participate in the collaborative process of producing the final manuscript. All analysts will sign a consent (ethics) document (LINK).

We identified our minimum number of analyses per data set by considering the number of effects needed in a meta-analysis to generate an estimate of heterogeneity (τ^2) with a 95% confidence interval that does not encompass zero. This minimum sample size is invariant regardless of τ^2 . This is because the same t-statistic value will be obtained by the same sample size regardless of variance (τ^2). We see this by first examining the formula for the standard error, SE for variance, (τ^2) or $SE(\tau^2)$ assuming normality in an underlying distribution of effect sizes (Knight 2000):

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

and then rearranging the above formula to show how the t-statistic is independent of τ^2 , as seen below.

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

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

Step 2: Primary Data Analyses

Analysis teams will register and answer a demographic and expertise survey (LINK). The survey collects information on the analysts current position and self-estimated breadth

and level of statistical expertise. We will then provide teams with the acoustic data set and request that they answer the following research question:

Do speakers acoustically modify utterances to signal atypical word combinations?

Once their analysis is complete, they will answer a structured survey (LINK), providing analysis technique, explanations of their analytical choices, quantitative results, and a statement describing their conclusions. They will also upload their analysis files (including the additionally derived data and text files that were used to extract and preprocess the acoustic data), their analysis code (if applicable), and a detailed journal-ready statistical methods section.

Step 3: Peer Reviews of Analyses

At a minimum, each analysis will be evaluated by four different reviewers, and each volunteer peer-reviewer will be randomly assigned to methods sections from at least four analyst teams (the exact number will depend on the number of analysis teams and peer reviewers recruited). Each peer reviewer will register and answer a demographic and expertise survey identical to that asked of the analysts. Reviewers will evaluate the methods of each of their assigned analyses one at a time in a sequence determined by the initiating authors. The sequences will be systematically assigned so that, if possible, each analysis is allocated to each position in the sequence for at least one reviewer. For instance, if each reviewer is assigned four analyses to review, then each analysis will be the first analysis assigned to at least one reviewer, the second analysis assigned to another reviewer, the third analysis assigned to yet another reviewer, and the fourth analysis assigned to a fourth reviewer. Balancing the order in which reviewers see the analyses controls for order effects, e.g. a reviewer might be less critical of the first methods section they read than the last. The process for a single reviewer will be as follows. First, the reviewer will receive a description of the methods of a single analysis. This will include the narrative methods section, the analysis

team's answers to our survey questions regarding their methods, including analysis code, and the data set. The reviewer will then be asked, in an online survey (LINK), to rate both the acoustic analysis and the statistical analysis on a scale of 0-100 based on these prompts:

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

The initiating authors (SC, JC, TR) will conduct the analyses outlined in this section. We will describe the variation in model specifications in several ways:

First, we will calculate summary statistics describing variation among analysis, including the nature and number of acoustic measures (e.g. f0 or duration), the operationalization and the temporal domain of measurement (e.g. mean of an interval or value at specified point in time), the nature and number of model parameters for both fixed and random effects [if applicable], the nature and reasoning behind inferential assessments (e.g. dichotomous decision based on p-values, ordinal decision based on Bayes factor), as well as the mean, standard deviation and range of effect sizes reported. We anticipate that the majority of statistical analyses will be expressible as a (generalized) linear regression model.

ADD FORMULA

Since teams will likely use outcome variable that substantially differ in their scales, we will standardize all reported effects by refitting all models with scaled outcome variables (the observed values subtracted from the mean divided by the standard deviation).

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. 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 measures such as Cohen’s d are inappropriate. Effect sizes will be defined as the estimate(s) of the critical predictor(s) (i.e. critical according to the analysis teams’ self-reported inferential criteria) divided by the standard error for the estimate(s).

Upon extracting the standardized effect sizes and standard errors for each analysis, the initiating authors will then fit a cross-classified Bayesian meta-analysis on the analyst team data using the multilevel regression model described below:

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

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

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

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

Effect size (δ_t) will be the outcome variable. The likelihood of the outcome variable is assumed to be normally distributed. Analysis teams will be included as a group-level effect (i.e., random effect). For all population-level parameters, the model will include regularizing, weakly informative priors (Gelman, 2017), which are normally distributed and centered at 0 with a standard deviation of 1.

A cauchy prior set at 0 with scale 1 will be used for τ . We will fit the model with 4000 iterations (2000 warm-up) and Hamiltonian Monte-Carlo sampling of the posterior distribution is carried out using 4 chains distributed across 4 processing cores. The analysis

will be conducted in R (R core team, 2020) and fit using `stan` (Stan, 2019) via the R package `brms` (Bürkner, 2019). The code for the model can be found here: INSERT LINK. We will quantify the extent to which the meta-analytic estimate is modulated by the following main predictors: The peer ratings of each analysis (numeric, 1-100) ... to be filled.

As a second step, we will explore the extent to which deviations from the meta-analytic mean by individual effect sizes relate to a series of predictors (see below): The deviation score, which serves as the dependent variable in this analysis, will be the

These analyses are secondary to our estimation of variation in effect sizes described above. We wish to quantify relationships among variables, but we have no a priori expectation of effect size and we will not make dichotomous decisions about statistical significance.

The following predictors will be used:

First, we include a measure of the ‘uniqueness’ of individual analyses for - the set of predictor parameters, - the set of random effect parameters, - the acoustic measurement.

The measure of the uniqueness of the set of model parameters is assessed by the Sorensen’s Similarity Index (SSI). The SSI is an index typically used in ecology research to compare species composition across sites. For our purposes, we will treat variables as species and individual analyses as sites. In order to generate an SSI for each analysis team, we will calculate the average of all pairwise Sorensen’s values for all pairs of analyses using the `betapart` package (Baselga et al. 2018) in R. We achieve this using the following formula:

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

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.

Second, we include a measure of conservativeness of the model specification defined by the number of random effect parameters.

Third, we include the self-proclaimed number of post-hoc changes to teams' acoustic measurements and the self-estimated number of models that they ran prior to settling on their final model.

Forth, we include the major acoustic dimension that has been measured to answer the research question. We will categorize each analysis into the following possible major acoustic dimensions: duration, amplitude, fundamental frequency, spectral properties

Fifth, we include the temporal window that the measurement is taken over defined by the target linguistic unit. We assume the following relevant linguistic units: segment, syllable, word, phrase.

Sixth, we include the following demographic factors about both the analysis teams: -

DESCRIBE HOW WE WILL LOOK AT THESE THINGS

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

427 interpreted as an answer to our primary question.

428 We wish to quantify relationships among variables, but we have no a priori expectation
429 of effect size and we will not make dichotomous decisions (such as statistical significance).

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

431 Analysts and initiating authors will discuss the limitations, results, and implications of
432 the study and collaborate on writing the final manuscript for review as a stage-2 Registered
433 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>