# Sample Size Matters in Calculating Pillai Scores

Joseph A. Stanley,<sup>1</sup> and Betsy Sneller<sup>2</sup>

<sup>1</sup> Linguistics Department, Brigham Young University, Provo, Utah 84602, United States

<sup>2</sup> Department of Linguistics, Languages, and Cultures, Michigan State University, East Lansing, Michigan 48824,

United States

## Abstract

Since their introduction to sociolinguistics by Hay et al. (2006), Pillai scores have become a standard metric for quantifying vowel overlap. However, there is no established threshold value for determining whether two vowels are merged, leading to conflicting ad hoc measures. Furthermore, as a parametric measure, Pillai scores are sensitive to sample size. In this paper, we use generated data from a simulated pair of underlyingly merged vowels to demonstrate (1) larger sample sizes yield reliably more accurate Pillai scores, (2) unequal group sizes across the two vowel classes is irrelevant in the calculation of Pillai scores, and (3) it takes many more data than many sociolinguistic studies typically analyze to return a reliably low Pillai score for underlyingly merged data. We provide some recommendations for maximizing reliability in the use of Pillai scores, and provide a formula to assist researchers in determining a reasonable threshold to use as an indicator of merged status given their sample size. We demonstrate these recommendations in action with a case study.

# I. INTRODUCTION

2 Quantifying vowel overlap is an important component of many linguistic studies, ranging from 3 sociolinguistics to laboratory phonology. While various methods of quantifying vowel overlap have 4 been proposed (see, e.g., Nycz & Hall-Lew 2013), Pillai scores have emerged in recent years as the 5 most commonly used method, particularly within sociolinguistics (Hay, Warren & Drager 2006; Nvcz 6 & Hall-Lew 2013), because of its ability to measure a distinction in multivariate space while also 7 accounting for fixed effects like phonological context. However, there is no standard value of Pillai 8 score that is broadly accepted to be a threshold for "merged" or "distinct"; indeed, individual studies 9 make this determination primarily by comparison between individuals in a single data set, to show that 10 some speakers are "more merged" while others are "less merged". In this paper, we provide a critical 11 look into how Pillai scores are calculated and reported.

12 We focus on demonstrating how sample size plays a major role in the resulting Pillai score. 13 We also show that the common approach of reporting Pillai scores alone is incomplete without also 14 reporting sample sizes and *p*-values. Using simulation data drawn from an underlyingly merged data 15 set, we show how larger samples produce lower Pillai scores (in other words, a more "merged" score). 16 We further demonstrate that within Pillai, it is the total *n* across both samples that matters, and provide 17 a formula that researchers can use to determine a threshold for "merged" given their own sample size. 18 We highlight some important takeaways about using Pillai scores to measure vowel overlap and the 19 potential risks for across- and within-study comparisons of speakers. We end with a case study which 20 demonstrates how researchers can implement the recommendations in this paper when analyzing real 21 data from sociolinguistic interviews.

## A. Pillai as a measure of vowel overlap

### 24 1. Pillai score overview

25 A multivariate analysis of variance (MANOVA) is an extension of the (univariate) analysis of 26 variance (ANOVA). The difference is that while an ANOVA evaluates whether the difference 27 between two or more groups in a single numeric variable can be predicted by some number of 28 categorical independent variables (such as the F2 of /u/across older and younger participants), a 29 MANOVA can evaluate two or more dependent numeric variables simultaneously (such as F1, F2, 30 F3, and duration of /u/ by older and younger participants). So, a researcher analyzing American 31 English vowels (which typically are differentiated primarily on vowel quality in F1-F2 space) could 32 use a MANOVA to see whether a speaker pronounces two historically distinct vowel classes 33 differently, while including the effects of duration and place of articulation as independent variables. 34 In this case, F1 and F2 would be the dependent variables, and vowel class, duration, and place of 35 articulation would be the independent variables.

36 Simplified somewhat, the null hypothesis of a MANOVA is that category membership (for 37 instance, two historically distinct vowel classes) offers no explanatory power for any of the 38 dependent variables. In the case of a MANOVA that is fit to vowel data, the null hypothesis is that 39 the two vowels are merged. In other words, there would be no way to guess which historic vowel 40 class a particular token came from by its acoustic measurements alone. Typically, the researcher's 41 aim is to find evidence to reject that hypothesis. We note that high p-values associated with 42 MANOVAs only indicate a lack of evidence to reject the null hypothesis that the two vowel classes 43 are the same, rather than evidence for the null hypothesis. MANOVA cannot prove that two vowels 44 are merged, only that there is little evidence to suggest that they're distinct.

45 There are four main test statistics associated with a MANOVA to compare what the data46 shows to the null hypothesis: Wilk's lambda, the Lawley-Hotelling trace, Roy's largest root, and the

47 Pillai-Bartlett trace. Of these four, the lattermost is the most robust for non-normally distributed data 48 and other violations of the assumptions of a MANOVA for tests that compare more than two groups 49 (Olson 1976). In sociophonetic data, where it is more common to compare only two groups, the Pillai-50 Bartlett trace has less advantage over the other three in statistical validity, but does benefit from being 51 both easy to run in commonly used statistical environments and relatively easy to interpret. For 52 introductory overviews of these four test statistics, including their mathematical definitions and 53 conceptual explanations, see Bray & Maxwell (1985: 27-29), John & Wichern (2012: 336), Rencher & 54 Christensen (2012: 169-188), and Upton & Cook (2014, "multivariate analysis of variance 55 (MANOVA)").

56 The Pillai-Bartlett trace, often called the Pillai score or occasionally just Pillai in linguistics 57 studies (a convention that we adopt here), ultimately comes from Pillai (1955) and Bartlett (1939). In 58 the simplest model, which predicts two dependent variables (e.g., F1 and F2) using a single two-level categorical variable (e.g., /a/and /3/), it returns a value that ranges between 0 and 1, with smaller 59 60 values occurring when there is greater overlap between the two groups in multivariate space, and larger 61 numbers for less overlap. In other words, small Pillai scores suggest a vowel merger. In reality, 62 determining whether a merger is present is not quite as simple as merely observing overlap. Two 63 vowels may occupy the same F1-F2 space but the distinction between phonemes may be maintained 64 though some other cue like voice quality (Di Paolo & Faber 1990), duration (Labov & Baranowski 65 2006), or vowel trajectory (Stanley 2020). We adopt a simplified approach here since our focus is to 66 explore the effects of sample size, but we acknowledge that there is more to merger than overlapping 67 midpoints in the F1-F2 space.

68

# 2. Meta-analyses of Pillai scores vs. other metrics

69 Prior to the introduction of Pillai scores to sociophonetics, perhaps the most common way to70 assess merger was through auditory coding. In its most basic form, the phonetically-trained researcher

71 would listen and evaluate whether there was a difference between the two sounds. But, as shown 72 below, Pillai scores are most commonly used on low vowels, which are notoriously difficult for 73 fieldworkers to accurately transcribe (Johnson 2010: 28–29). In fact, Moulton (1968: 464) rather 74 strongly states that early fieldworkers for the Linguistic Atlas Projects were "hopelessly and humanly 75 incompetent at transcribing phonetically the low and low back vowels that they heard from their 76 informants." Fortunately, formulaically quantifying overlap provides a less subjective measure for 77 vowel merger.

78 While Pillai scores are currently the most common metric for quantifying mergers, especially 79 within sociolinguistic work, there are also other approaches available. In addition to the auditory 80 coding mentioned above, much early work used the Euclidean Distance between the point 81 representing the mean F1 and mean F2 of one vowel class and the point representing the mean F1 82 and mean F2 of a second vowel class as the primary metric for vowel merger (Hay, Warren & Drager 83 2006; Nycz & Hall-Lew 2013; Han & Kang 2013; Hall-Lew 2013). This early approach is not particularly satisfying, as it fails to take into account the distributional properties of the data, including 84 85 the degree of overlap and the distribution of tokens within a vowel class (Kelley & Tucker 2020: 137). 86 The Spectral Overlap Assessment Metric (SOAM; Wassink 2006), which calculates overlap between 87 ellipses or ellipsoids fitted to the vowel distribution (see Wassink 2006 for details on the fitting) and 88 calculates the area or volume of their overlap, is one method adopted and recommended in other 89 sociophonetic work (Di Paolo, Yaeger-Dror & Wassink 2011: 103; Kendall & Fridland 2021: 56). We 90 refer interested readers to Nycz and Hall-Lew (2013) and Kelley and Tucker (2020) for in-depth 91 assessments of these measures and several others, compared with Pillai scores. Kelley and Tucker 92 assess four different metrics of vowel overlap, using a Monte Carlo simulation to test accuracy, and 93 find that Pillai scores produced the most accurate and precise values when compared to ground truth

94 values in their simulated data. They also recommend Pillai scores when sample sizes are "small," which95 they define as 30 observations per group.

96 For sociophonetic data, and especially for naturalistic sociophonetic data, obtaining more than 97 30 observations per group is not always feasible, making Pillai an especially valuable tool for 98 sociophonetics. As perhaps foregrounded by the title of this current paper, sample size plays an 99 important role in the resulting Pillai score, which we highlight in detail in Sections 2-3. We note, 100 however, that this is generally true of all measures of vowel overlap compared in the overview papers 101 mentioned above (Nycz and Hall Lew 2013; Kelley and Tucker 2020). Indeed, any measure that uses 102 means, standard deviations, and/or variance as inputs will be impacted by sample size. Likewise, larger 103 sample sizes impact statistical significance, with larger sample sizes leading to smaller p-values. Our 104 aim here is not to simply recommend that researchers obtain larger sample sizes, which in many cases 105 is either not possible or may be at odds with other important data collection considerations, but rather 106 to elucidate *how* sample size impacts resulting scores, so that researchers can be best informed when 107 using Pillai as a measure of vowel overlap.

108 Perhaps an even bigger concern than total sample size with naturalistic data is the near 109 impossibility of obtaining balanced token counts across two categories from naturalistic data. In 110 wordlist data, researchers can carefully craft their wordlist to obtain balanced data: this typically means 111 obtaining the same number of observations in each vowel class and ensuring that the wordlist is 112 balanced for additional factors, such as phonological context. In naturalistic (i.e., conversational) 113 speech, there is no way to ensure that a speaker produces balanced token counts across vowel category 114 and across phonological contexts. As a result, researchers investigating something like vowel overlap 115 need to understand precisely how Pillai scores may be affected by unbalanced data.

Finally, while Pillai has emerged as the best option so far for most sociophonetic data, thereis one additional concern to address, which is that MANOVAs assume the data within each group

118 follows a multivariate normal distribution (Bray & Maxwell 1985). Formant measurements for a given 119 vowel from a given speaker in naturalistic vowel data, being nonnormally distributed, do not fit this 120 assumption (though see Whalen & Chen 2019 for some evidence that vowel formant data, even with 121 coarticulation, can be normally distributed). Recent work has suggested Bhattacharyya's Affinity<sup>i</sup> as a 122 better measure for non-normally distributed data (Fieberg & Kochanny 2005; Johnson 2015), and has 123 been taken up by some subsequent work (Strelluf 2018; Warren 2018; Jensen & Braber 2021). While 124 the benefits of Bhattacharyya's Affinity being a nonparametric measure make it particularly appealing 125 for the kind of distributions of data found in naturalistic speech, there is not yet a mechanism for easily incorporating fixed effects like phonological context or speaker age." Furthermore, 126 127 Bhattacharyya's Affinity also works best on a relatively large data set of over 30 observations per 128 category (Seaman et al. 1999). While future work may make nonparametric methods like 129 Bhattacharyya's Affinity more easily integratable with the kind of statistical models linguists typically 130 use, for now we focus on Pillai score and the considerations necessary to make its use maximally 131 standardized.

132

133

# 3. How Pillai scores are used to measure overlap in sociophonetics

134 Pillai scores have been used to analyze a variety of phenomena in several languages. Perhaps 135 the most common application of Pillai scores is to measure overlap between  $/\alpha/and /3/in$  North 136 American English (Hall-Lew 2013; Kendall & Fridland 2017; Havenhill 2015; Stanford et al. 2019 137 inter alia). In fact, using Pillai scores to measure this merger was explicitly recommended in Becker 138 (2019), a volume of different studies all analyzing the spread of the  $/\alpha$ - $\sigma$ / merger in North American 139 English. Many conditioned mergers in English have been quantified with Pillai scores as well (Schmidt, 140 Diskin-Holdaway & Loakes 2021; Austin 2020; Freeman 2021; Newbert 2021 inter alia). To a lesser 141 extent, Pillai scores have been used to analyze vowels that are marginally contrastive (Galician:

Amengual & Chamorro 2015; Italian: Nadeu & Renwick 2016; Bangla: Islam & Ahmed 2020;
Hawai'ian: Kettig 2021; Swiss German: Joo, Schwarz & Page 2018; Austrian German: Sloos 2013).
Some more innovative uses of Pillai scores include using them to analyze tones in varieties of
Cantonese (Fung & Lee 2019; Tse 2018) and in Spanish fricative mergers (Regan 2020).

146 Pillai scores have also been used to quantify splits, chiefly among phonological low vowels. 147 Fisher et al. (2015) assess the degree to which the Philadelphia short-a split is found in their sample 148 of Philadelphians. Relatedly, Hall-Lew et al. (2017) look at the BATH-TRAP split in Scottish Parliament 149 data. Brozovsky (2020) uses Pillai scores to measure the raising (and separation) of prenasal  $/\alpha$ , using 150 Pillai scores to measure the overlap between prenasal  $/\alpha$  and preobstruent  $/\alpha$  in Taiwanese Texans. 151 In a study looking at the possible effect of salience on a given lexical item compared to the rest of its 152 canonical vowel class, Bray (2021) analyzes the lowered realization of the vowel in *bockey* compared to 153 other relatively more raised /a/ tokens in professional American hockey players.

As highlighted by the selection of citations above, since being introduced to the field, Pillai scores have become widespread in sociophonetic studies. With the support of meta-analyses that compare other competing measures, Pillai has become a useful go-to for measuring both the overlap and the distinction of speakers' pronunciation of two phonological categories, especially in sociolinguistic studies that need to compare across individual speakers. While Pillai scores are clearly a valuable tool in measuring vowel overlap, there remain some outstanding issues with using it. In the following sections, we highlight some of these issues.

161

- 162 B. Issues with Pillai scores
- 163 *4. What is considered merged?*

As useful as Pillai scores are for quantifying the degree of overlap, they do not necessarilyanswer researchers' underlying question of whether two vowels are merged. Pillai scores range from

166 0 to 1, but there is no agreed-upon cutoff value or threshold for determining whether the two groups 167 are merged or not. As a result, many studies rely on an ad-hoc threshold to interpret the merged status 168 of their speakers. Some work has suggested specific thresholds for mergers. For instance, Jibson 169 (2021) suggested a Pillai threshold of 0.3 as an indicator of "merged" status, after a shuffling procedure identified 0.3 as the 95th percentile of "merged" between 20 tokens of two vowel classes from his 170 171 speakers. Wassink (2006) likewise suggests some provisional thresholds for SOAM, where 0-20% overlap represents "distinct", 20-40% represents "partially merged", and >40% represents "merged". 172 Relying on provisional or ad-hoc thresholds, however, is risky because sample size is likely not 173 174 comparable across studies or even between speakers.

175 One solution for determining whether a given Pillai score should be interpreted as an 176 indication of "merger" is to examine the *p*-values that are associated with the MANOVA model from 177 which the Pillai scores are generated. The model assumes that the vowel variable contributes no 178 information to differentiating between two groups. In other words, it assumes the two vowels are 179 underlyingly merged. A small p-value associated with the vowel variable would provide evidence 180 against that null hypothesis, allowing the researcher to conclude that the difference between the two 181 groups is likely true (i.e., that the speaker is not merged). We note that Pillai scores and *p*-values are 182 inversely correlated: lower Pillai scores typically accompany higher *p*-values. In fact, since Pillai scores 183 are just test statistics, they and *p*-values are functions of each other. Nevertheless, *p*-values are not 184 typically reported in sociophonetic studies that use Pillai (some exceptions include Wong & Hall-Lew 185 2014; Nadeu & Renwick 2016; Amengual & Chamorro 2015; Berry 2018; Sloos 2013).

186 There are a few points of caution to make about using and interpreting *p*-values, as we discuss 187 throughout this paper. One of these concerns the potential distinction between a statistically 188 significant difference (as defined by the model) and a ground truth difference for speakers. For a 189 speech community that has a ground truth merger in two vowel categories, there will be no difference 190 in their perception of these two vowel categories. However, as sample size increases, so does the 191 likelihood of a model returning a p-value below a given significance threshold (typically 0.05); it is 192 possible that even for pairs of sounds that are truly merged, a sufficiently large dataset can interpret 193 random variance in the data as a meaningful difference, as shown in the experiments in Section 3. An 194 additional caution to make regarding interpreting *p*-values alone as an indicator of merger is that a 195 statistically significant difference in formant values may not map onto a perceptible difference for the 196 human auditory system (see, e.g., Kewley-Port and Watson, 1994). P-values alone can likewise be 197 misleading in the opposite direction: previous work has shown that speakers and listeners can produce 198 and perceive reliable but small differences, including sub-phonemic differences, in what may otherwise 199 appear to be merged sounds (for instance, with cases of incomplete neutralization; Warner et al. 2004; 200 Pfiffner 2021). Vowel distinctions that are maintained by small effect sizes, or by sub-phonemic 201 distinctions not captured by the measurements in the model, may appear artificially to be merged 202 according to a *p*-value because it takes more data for smaller differences to be detected by the statistical 203 model. For these reasons, additional information such as Pillai scores can aid in the interpretation of 204 *p*-values, and vice versa.

# 205 *5. Sample size*

As suggested by the parametric nature of Pillai, and the discussion of ad-hoc thresholds above, a major component of deciding which threshold should be used to determine merger status is the number of tokens being analyzed. Previous work on Pillai scores and sample size have expressed concern over too-small sample sizes (Gorman & Johnson 2013), and over unbalanced sample sizes across the two vowel classes being analyzed (Nycz & Hall-Lew 2013; Johnson 2015).

Despite sample size having a major impact on Pillai scores, most studies in sociophonetics
that use Pillai as a measure of vowel overlap do not also clearly report sample size (exceptions include
Wong & Hall-Lew 2014; Holland & Brandenburg 2017; Berry 2018; Berry & Ernestus 2018). This in

turn makes it difficult both to assess the findings in an individual paper and to compare speakers within and across studies. In the simulation experiments below, we show just how important sample size is to the resulting Pillai score, and provide a formula to calculate a recommended Pillai score threshold for merger status, given a particular sample size.

218

# II. METHODS

In this section, we present the results of Monte Carlo simulation experiments designed to test the effect of different sample sizes on resulting Pillai score. In these simulations, we create two vowel classes that are perfectly merged underlyingly, and alter (1) the sample sizes between the two vowel classes to test the effect of unbalanced samples across vowel classes, (2) the overall sample size, considering both vowel classes together, to test the effect of unbalanced total sample size across speakers, and (3) the correlation between the simulated F1 and F2 formant frequencies, to test whether correlations (like those typically found in naturalistic vowel data) influence the results.

226

# C. Data generation

227 We generated data using a Monte Carlo simulation (Metropolis & Ulam 1949). This is a 228 procedure where random draws are taken from an underlying probability distribution or existing 229 dataset and analyzed. This process is repeated independently many times and the information about 230 each iteration is aggregated. To begin the simulation, a bivariate normal distribution<sup>iii</sup> was generated 231 in R to simulate a single theoretical underlying vowel in the F1-F2 dimension for a single theoretical 232 speaker. For the sake of simplicity, the mean for F1 and F2 were both set to zero and the standard 233 deviation was 1. In Experiments 1 and 2, the correlation coefficient between the formants was zero, 234 producing a circular (rather than elliptical) distribution. In naturalistic speech, however, vowel data 235 typically has some degree of correlation between F1 and F2, resulting in elliptical distributions, so in 236 Experiment 3, we manipulated the correlation coefficient between F1 and F2, to test whether the 237 results from Experiments 1 and 2 hold for data with different degrees of correlation. To simulate 238 vowel data, random draws were taken from that single bivariate normal distribution, and assigned to 239 one of two "vowel class" labels. We note that it is somewhat nonsensical to refer to these generated 240 numbers as "vowels", especially since Pillai scores can be calculated on non-vowel data. However, 241 since the majority of Pillai scores in sociophonetic analyses are based on vowels, for clarity, we will 242 refer to this simulated data points as "vowels" and their arbitrary groups as "vowel classes." These 243 random draws represent a linguist sampling data from our theoretical speaker. For the sake of 244 illustration, we will say 30 such observations were generated. These 30 observations were treated as tokens from a single underlying vowel class.<sup>iv</sup> Another 30 random draws were then taken from the 245 246 same bivariate normal distribution. These 30 observations were treated as tokens from a different 247 underlying vowel class. Generating two groups from the same underlying distribution therefore creates 248 a simulated pair of merged vowels. In theory, the two simulated vowel classes should not be statistically 249 different from each other in any way because they were drawn from the same underlying distribution.

250

## 251 D. Three experiments

For this study, we ran three experiments. In Experiment 1, the two simulated vowel classes for each "speaker" were of equal size. We began with a sample size of 5 observations per vowel class. We then moved on to 6 observations per vowel class, and so on, until we reached 100 observations per vowel class. For each of these 96 sample sizes, we repeated the simulation 1000 times, each representing a different instance of a linguist sampling data from that one underlyingly merged speaker. This produced 96,000 pairs of simulated vowel data, where each pair consisted of equal-sized vowel classes, enabling us to test the effect of overall sample size on resulting Pillai score.

In Experiment 2, we varied the sample size between the two vowel classes for each "speaker".We began with 5 tokens from one vowel and 6 from another. We then took 5 tokens of one and 7

from the other. We increased the size of the second group by steps of 1 until it contained 100 tokens.
We then repeated this process with the first group having 6 tokens, and increased the second group
from 5 to 100 in steps of 1. We iterated over these steps, increasing the sample size of the first group
up to 100, thereby generating pairs of vowel data where every combination of sample sizes from 5 to
100 was represented. We repeated this simulation 100 times per combination. This produced 921,600
pairs of simulated vowel data, where each pair consisted of different-sized vowel classes, enabling us
to test the effect of unbalanced vowel class size on resulting Pillai score.

268 In Experiment 3, we then varied the correlation between the simulated formant data, 269 modifying the shape of the underlying bivariate normal distribution from circular to elliptical using 270 the mvrnorm() function in the MASS package (Venables, Ripley & Venables 2002). As with 271 Experiments 1 and 2, this produced an underlyingly merged pool, which we could then sample from 272 to generate our two "vowel classes". Following the methods described in Experiment 2, we generated 273 data with sample sizes ranging from 5 to 100 tokens per vowel, though for the sake of processing 274 time, we only chose samples that were multiples of 5. For each combination of sample sizes then, we 275 generated datasets with correlation coefficients ranging from 0 to 0.9, in intervals of 0.1. 100 vowel 276 pairs per combination of sample sizes and correlation coefficients were produced, resulting in 361,000 277 new sets of simulated vowel data, enabling us to test the effect of correlation on Pillai scores, in 278 conjunction with sample sizes and unbalanced vowel class sizes.

Across all experiments, Pillai scores were calculated for each pair of simulated vowel data. Pillai scores were calculated by fitting a MANOVA model to the data using the manova() function in R. The simulated F1-F2 measurements were the dependent variables and the vowel class was the only independent variable. While it would be possible to incorporate additional simulated independent variables such as place of articulation and duration into the MANOVA, we consider this to be beyond 284 the scope of the current paper, which focuses on the effect of sample size on resulting Pillai scores. 285 We therefore include only historical vowel class in our models, and leave more complex MANOVA 286 models to future work. The Pillai scores and p-values associated with the vowel class variable were 287 then extracted from that MANOVA model. To reiterate, the Pillai scores for all of these distributions 288 should be very close to zero (indicating complete overlap) because every vowel pair was generated 289 from the same underlying bivariate normal distribution. Because the data is randomly generated, some 290 Pillai scores will be higher than others, but by rerunning the simulation many times per sample size, 291 we can begin to see patterns that may emerge at a given sample size.

The coding for this study was done in the R programming language (R Core Team 2021) with the help of the tidyverse suite of packages (Wickham et al. 2019) and joeyr (Stanley 2021).

Visualizations were generated using ggplot2 (Wickham 2015) and see (Lüdecke, Patil, et al. 2021) with

color palettes from ggthemes (Arnold 2018) and scico (Pedersen & Crameri 2020).

#### 296 III. RESULTS

297

# E. Experiment 1: Equal sample sizes

298 To address how sample size affects Pillai scores, we first present the results from Experiment 299 1, where the two simulated vowel classes were the same size. Before inspecting the results of all sample 300 sizes though, it is important to understand how the 1000 Pillai scores were distributed within a given 301 sample size. Figure 1 shows two different views of the distribution of Pillai scores when the sample 302 size for both groups was 10. We see that the distribution of points representing resulting Pillai scores 303 is rather wide, a consequence of using such a small sample size for inferential statistics, ranging from 304 less than 0.001 to 0.568. Much of the data is clustered near the bottom of the plot but there is a long 305 "tail" extending upwards. This is not a haphazard pattern, but rather follows a distribution that can

be transformed into an *F* distribution and reflects the underlying mathematical properties of how Pillai
scores are calculated (cf Rencher & Christensen 2012: 182). For this particular sample size, the mean
Pillai score was 0.104, the median was 0.077, and the 95<sup>th</sup> percentile was 0.294. As seen below, these
numbers change depending on the sample size, but the underlying distribution of the Pillai scores is
consistent across sample sizes. We show this distribution to dispel any misconceptions that Pillai
scores are uniformly distributed within a particular range, and to highlight that generally they fall near
the lower end of the distribution.



**315** Figure 1: Distribution of Pillai scores on 1000 pairs of simulated groups, each with a size of

10.

317 With that distribution in mind, we can now zoom out to view all samples at once. Figure 2 318 shows all 96,000 Pillai scores by their sample size. Though present at all sample sizes, the "bottom-319 heavy" distribution shown in in Figure 1 is not displayed in Figure 2, in order to make the general 320 trend across samples easier to see. It is immediately apparent that larger groups more consistently 321 produced lower Pillai scores. With very small sample sizes (groups of fewer than 10 observations per 322 vowel class), Pillai scores were quite high. For these small sample sizes, Pillai scores were sometimes 323 closer to 1 than 0, even for these underlyingly merged vowel classes which should, in principle, return 324 a Pillai score of 0. In other words, these small sample sizes sometimes resulted in very misleading Pillai 325 scores that may cause a researcher to interpret two vowel classes as distinct even if the true underlying 326 distribution was perfectly merged. As the sample size increases, Pillai scores were more consistently 327 low, as we would expect for underlyingly merged vowels.

328 The black line overlayed on Figure 2 indicates the 95<sup>th</sup> percentile for each sample size. This 329 line also very closely corresponds to the threshold for vowel class being statistically significant in the 330 MANOVA models: almost all points above that line had p-values less than 0.05 while almost all points 331 below it had greater *p*-values. Because we are modeling the null hypothesis (*i.e.* underlyingly merged 332 vowel classes), the distribution of *p*-values is uniform. It therefore is unsurprising, and in fact expected, 333 that the highest 5% of Pillai scores within a given sample size also return p-values less than 0.05. This 334 line shows that if there are just 10 observations per vowel class, 95% of the Pillai scores were under 335 0.3. However, as is evident in Figure 2, a threshold of 0.3 is only applicable to a sample size of 10 per group since Pillai scores decrease with larger samples. For example, the 95th percentile of returned 336 337 Pillai scores does not drop to 0.1 until there are 30 observations per group.



- 340 sample size.
- 341

# F. Experiment 2: Unequal sample sizes

While the previous section found that sample size affects Pillai scores in a predictable way, with larger samples returning more reliable Pillai scores, in this section we conduct further simulations to explore what effect, if any, an unbalanced sample has on Pillai scores. Unless data collection is carefully controlled to include a fixed number of tokens per vowel class, Pillai scores are run on vowel classes that are not comprised of the same number of tokens. Here we ask: what effect do grossly unbalanced groups have on Pillai scores?



**350** Figure 3 (color online): Mean Pillai scores for all simulations.

Figure 3 shows the mean Pillai scores from the 100 simulations for each combination of sample sizes between the two vowel classes. Going from the bottom left corner (two small sample sizes from each vowel class) to the top right corner (two large sample sizes from each vowel class), we see a general trend of decreasing Pillai scores as sample sizes increase. Reflecting Figure 2, we see these decreasing Pillai scores dropping more sharply as sample sizes are small (under around 30 tokens per vowel class).

358 A surprising pattern emerges when we look closely at the resulting Pillai scores from unequal 359 sample sizes: namely, that unequal samples across the two vowel classes do not impact Pillai score— 360 it is only the *total* sample size taking both vowel classes together that matters. For example, the mean 361 Pillai score for a pair of 10 tokens and 50 tokens drawn from this underlyingly merged distribution is 362 0.0358 and the mean Pillai score for a pair of 20 and 40 tokens is 0.0355. A pair of 30 tokens and 30 363 tokens drawn from an underlyingly merged distribution is nearly identical: 0.0381. We see this pattern 364 visually reflected along the diagonal between the top left corner and the bottom right corner of Figure 365 3, which shows a symmetrical resulting Pillai score for all unequal pairs of samples that sum to the 366 same total. In other words, we find that neither the existence of an unequal sample size between vowel 367 classes nor the degree of unequalness impact Pillai score.

These findings bring us to recommend simply using as many tokens as researchers have available for an individual speaker, to bring the *total* sample size as high as possible. As seen in Section 5 below, when comparing across speakers, it is important to recall that total sample size impacts the resulting Pillai score, and we recommend normalizing and reporting total sample sizes across speakers in a study to make resulting Pillai scores maximally comparable.

373

374

# G. Experiment 3: Correlated F1 and F2 values

One potential caveat for Experiments 1 and 2 is that they are based on uncorrelated F1 and F2 values (producing a circular distribution in F1-F2 space), while in naturalistic speech vowel formant data is typically correlated (elliptical in F1-F2 space). To evaluate whether the correlation of the vowel formants affects Pillai scores, our final experiment explicitly tests the effect that correlated dependent variables have on Pillai scores by generating bivariate normal distributions with various degrees of correlation. Like the previous two experiments, Pillai scores were calculated, only this time the effect of sample size, unequal group sizes, and correlation were explored. 382 Since this experiment modified three variables (group 1 size, group 2 size, and correlation), 383 visualization of all the data is less straightforward. Instead, we ran a linear regression model on these 384 Pillai scores, with the log-transformed Pillai score as the dependent variable and the log-transformed 385 total sample size and the correlation as predictors (Table 1). We find that correlation was not a significant predictor in these 361,000 simulations, meaning there was no significant difference between 386 387 the Pillai scores of uncorrelated data and correlated data, given a particular total sample size. We also find that unbalanced data did not affect Pillai scores, even in these correlated datasets. Thus we can 388 389 be confident that our results, and the implications and recommendations drawn from them, are 390 applicable even to correlated, unbalanced data like what is typically found in real vowel formant 391 measurements.

392

393 Table 1: Summary statistics of a linear regression model on Pillai score, showing that, despite 394 the very large number of observations, correlation is not a significant predictor.

	Estimate	Std. Error	<i>t</i> -value	<i>p</i> -value
Intercept	0.243484	0.023980	10.153	< 0.001
$\log(n)$		0.005114	-199.965	<0.001
	1.022667			
correlation	_	0.007372	-0.451	0.652
	0.003328			
Adjusted R <sup>2</sup>	= 0.09972; F(2, 360)	(997) = 19,993.16; p	< 0.001.	•

395

396

397

# 399 IV. IMPLICATIONS

400

### A. Choosing a threshold based on sample size

As a result of these simulations, we are in a position to recommend a formula that researchers can use as a guide to determine whether a given resulting Pillai score indicates a merger. We sought a formula that was a function of sample size and could provide the 95<sup>th</sup> percentile for Pillai scores given that two groups come from the same underlying distribution. In other words, we wanted to find the function for the black line in Figure 2. We chose the 95<sup>th</sup> percentile for this distribution to align with the common benchmark of p < 0.05.

407 The formula is based on the observation that by taking the natural log of the Pillai score and 408 the natural log of half of the total sample size of both groups (essentially log-transforming the axes in Figure 2), the 95<sup>th</sup> percentile per sample size followed a straight line with an intercept of 1 and a slope 409 410 of -1. Because "half of the total sample size" can be confusing, we begin by simplifying this to a 411 variable *m*, which represents n/2 where *n* is the total sample size across both groups. In other words, 412 *m* is the average sample size per group. These intercepts were determined by rerunning Experiment 1 413 100 more times and fitting a simple linear regression to the 95% confidence interval for the Pillai 414 scores in each iteration. We found that a line with an intercept of 1 and a slope of -1 was within the 415 95% confidence interval of those parameters approximately 95% of the time. We therefore assume 416 that these parameters are correct in the probability density function for the line we seek to model. The 417 formula therefore begins as pillai = 1 - m. But because this straight line only makes sense when both 418 x and y are log-transformed, it must be modified to be log(pillai) = 1 - log(m). At this point, we solve for *pillai*, using e<sup>x</sup> as the antilog function,  $p_{95} = e^{1 - log (m)}$ . Reducing view produces the formula 419 420 that we recommend for determining a Pillai score cut off,

421

$$p_{95} = \frac{e}{m} \tag{1}$$

424 where  $p_{95}$  is the 95<sup>th</sup> percentile of Pillai scores given the average sample size per group *m*. The 425 curve that is generated by Equation 1 very nearly follows the black line in Figure 2, which represents 426 the 95<sup>th</sup> percentile of Pillai scores given the average sample size (which is the same as the sample size 427 for each group in that figure since both groups were the same size). We believe this formula can be 428 fruitfully used to determine a reasonable cutoff for "merged" status for a given total sample size.

429 Using Equation 1, we see that, for example, a sample of 20 total observations would produce 430 a Pillai score less than or equal to 0.2718 95% of the time if the two vowel classes were underlyingly 431 merged (we note that this is a close approximation to the cutoff of 0.3 that Jibson (2021) chose for 432 his total sample size of 20). For illustration, Table 2 presents threshold suggestions drawn from this 433 formula for different total sample sizes, highlighting that it takes a great deal of data to reliably return 434 Pillai scores that are close to zero for underlyingly merged data (recall that Pillai scores closer to zero 435 reflect a more "merged" production). We recommend that researchers use Equation 1 in conjunction 436 with *p*-values to make a more informed decision about the merged status of two vowel classes rather 437 than an ad hoc or arbitrary cutoff.

- 438
- 439

Table 2: Pillai thresholds at various sample sizes based on Equation 1.

Total sample size	Pillai threshold for "merged"
1	8
20	0.2718
20	0.2710
40	0 1350
40	0.1557
60	0.0006
00	0.0900
	0.0700
80	0.0680
100	0.0544
100	

120	0.0453

441

# H. Sample size matters across speakers but not across vowels

442 One important takeaway from these simulations is that it takes a relatively large amount of 443 data to reliably (meaning 95% of the time) return a low Pillai score such as 0.1 from two underlyingly 444 merged vowels. In an analysis of English  $/\alpha/$  and /3/, for example, conversational data can typically 445 provide a sufficient number of observations for a robust analysis of overlap. However, few studies 446 that analyze wordlist data contain many more than 10 tokens of these two vowels. And even within 447 long-form conversational data, if the research question focuses on an infrequent phonological variable, 448 the total sample size quickly drops. Because total sample size has a major impact on the resulting Pillai 449 score, we recommend researchers choose a relevant threshold for "merger" status, based on their total 450 sample size, and use it in conjunction with the resulting *p*-value to make a determination about merger 451 status for their data.

452 Perhaps the most surprising takeaway from these simulations, for both authors, was that 453 although Pillai is not a nonparametric test, it does not actually matter if the token counts across the 454 two categories being investigated are unequal. Instead, the most critical consideration is the total number of tokens, summed across both categories. This is particularly important for naturalistic 455 456 sociolinguistic work, which relies on casual conversation rather than carefully constructed word lists 457 for data, meaning that is it often not possible to obtain balanced token counts across categories. 458 Following the results of Experiment 2, we can reassure researchers that unbalanced tokens across 459 vowel classes will not impact the resulting Pillai score. Instead, and following the results of Experiment 460 1, we recommend using as many total tokens possible for an analysis of a single speaker, regardless of 461 unbalanced samples across vowel classes for that speaker.

462 At the same time, any study aiming to compare the "merger" status of multiple individual 463 speakers should take into account the total sample sizes for each speaker, and especially consider the 464 fact that sample sizes (and therefore, the interpretation of resulting Pillai scores) may be different 465 across speakers. One of the primary goals for a robust measure of vowel overlap in sociolinguistics 466 has been to track the development of vowel mergers and splits across speakers in a given corpus (Nycz 467 & Hall-Lew 2013; Johnson 2010; Strelluf 2018; Labov et al. 2016) to analyze the trajectory of large-468 scale language change over time. Because distributions with lower token counts produce inflated Pillai 469 scores, this means speakers who are less talkative will artificially appear to have more distinct vowels 470 than speakers who are very talkative.<sup>vi</sup> We recommend using one of two ways to account for unequal 471 sample sizes across speakers.

472 One option is for researchers to conduct an analysis of individual speakers, incorporating all 473 the relevant pieces of evidence (the recommended Pillai threshold given an individual speaker's sample size, Pillai scores, p-values, and visualizations), to make a determination about the merged status of 474 475 each speaker. Section V.A presents an example of this method in action, where we demonstrate how 476 we leveraged all these pieces of evidence together to try to understand the merged status of each 477 speaker. This method allows researchers to obtain a fairly robust understanding of individuals as they 478 compare to each other and across styles. However, since this method requires researchers to synthesize 479 a number of gradient measures into discrete categories (such as "merged", "ummerged", and perhaps 480 "partially merged"), it makes it more difficult to track fine grained changes in the development of a 481 merger across a large speech community over time.

482 For researchers aiming to analyze fine grained differences over real or apparent time across
483 many speakers, however, it may be more beneficial to keep Pillai as a gradient measure. To analyze
484 Pillai across many speakers at once, we recommend analyzing the same number of tokens per speaker.

485 Our recommendations for how best to do this are discussed in detail in Section V.B, and an example486 of the R code used to apply this recommendation is provided in the supplementary file.

487 A similar issue arises when comparing merger status from a single individual speaker but across 488 multiple speech styles. In particular, sociophonetic work often compares casual speech styles like 489 conversation to more formal speech styles, such as a minimal pair list or a reading task. Critically for 490 Pillai scores, since lower token counts will artificially appear more distinct, speech styles such as word 491 lists and reading tasks with lower token counts will likewise artificially appear more distinct than speech 492 styles with higher token counts (such as casual speech). In sum, there is a strong risk that a Pillai score 493 difference between speech styles may be interpreted as a stylistic difference (whereby speakers appear 494 to be maintaining a distinction in wordlists that they do not maintain in casual speech), when the effect 495 is entirely driven by differences in total sample sizes across the two styles. It's not uncommon, for 496 instance, to hear of speakers apparently undoing mergers in read speech in comparison to their casual 497 speech (cf. Labov 1994: 80; Berry 2018; Berry & Ernestus 2018); while Labov (1994) correctly points 498 out that apparent unmergers in careful speech may be speakers hypercorrecting in response to orthographic differences across historically distinct vowel classes, we can add that distinctions 499 500 measured by Pillai score will additionally be impacted by sample size. We urge researchers who use 501 Pillai as a metric of overlap to pay close attention to differences in total sample sizes across speech style, v<sup>ü</sup> and either control for total sample size across styles or adjust their threshold of "merger" in 502 503 each style accordingly, following the same recommendations provided above and elaborated on in 504 Section V.

505

506

### I. What to report when using Pillai scores

507 Finally, we end with some general recommendations for what researchers should include in508 their results when using Pillai scores as a measure of overlap, in addition to the model specification

509 (i.e, the independent and dependent variables). First, because total sample size impacts Pillai score so 510 strongly, we recommend always reporting total sample size per speaker (or per style, for studies that 511 compare merger across speech styles). This practice will have the added benefit of enabling better 512 comparisons across sociolinguistic studies, in turn enabling researchers to gain a clearer picture of the 513 large-scale spread of some ongoing mergers such as the  $/\alpha/-/\mathfrak{0}/$  merger.

514 Second, in addition to reporting and controlling for total sample sizes, we recommend 515 reporting *p*-values where possible. In other words, where studies report individual speakers' (or styles') 516 Pillai scores, those should always include the total sample size and the *p*-value as well. We are not 517 necessarily advocating for an increased reliance upon p-values as a binary meaningful threshold, 518 particularly since some statisticians are urging quantitative researchers to abandon declarations of 519 "statistical significance" and to instead understand *p*-values as one gradient measure in concert with 520 additional evidence (Wasserstein, Schirm & Lazar 2019). However, reporting a Pillai score without its 521 accompanying p-value is akin to reporting a regression estimate without a p-value, a p-value without 522 an effect size, or a mean without a standard deviation. A Pillai score is a test statistic and should be 523 reported as such. These two numbers in conjunction paint a better picture of the merged status of a 524 pair of vowels than either one in isolation.

525

526 V. A CASE STUDY

The following is a case study to illustrate how one might conduct an analysis of vowel merger using the recommendations in this paper. As illustrative data, we draw from sociolinguistic interviews conducted and analyzed by the first author with residents of southwest Washington State, a region where the low back vowels (/ $\alpha$ - $\sigma$ /) are often merged, though with some indication of separation in some speakers (see Stanley 2020 for more details). On average across the 52 participants analyzed here, interviews lasted 46 minutes and yielded 143 tokens containing either  $/\alpha$  or  $/\mathfrak{o}/$  in preobstruent position. After the interviews, 30 of those participants then read a wordlist containing another 20 tokens in a more careful style.

- 535
- 536

# J. Analysis of individual speakers

We performed a MANOVA on each speaker's F1 and F2 measurements, separately for each style, with historic vowel class as the only predictor. Given the number of observations produced by each speaker, Equation 1 was implemented to establish a potential cutoff value for each style. The *p*values and the Pillai scores were extracted from the MANOVA model and the latter were compared to the cutoff values.

542 In this sample, there are some cases where the evidence overwhelmingly points towards a 543 merger. For example, 48-year-old Donna produced 179 low back vowels in the conversational portion 544 of her interview. With that many tokens, Equation 1 suggests that if her vowels were underlyingly 545 merged, her Pillai score should be less than 0.0304 about 95% of the time. In other words, anything 546 less than 0.0304 would be evidence for a merger. As it turns out, the MANOVA model performed on 547 her data yielded a Pillai score of 0.0289 with a p-value of 0.0756. The fact that her Pillai score is less 548 than the threshold for her token count and that the *p*-value above 0.05 means that historical vowel 549 class category does not predict Donna's acoustic realization, suggesting that Donna's two vowel 550 classes are likely underlyingly merged. In the wordlist data, there were only 20 tokens, so the suggested 551 cutoff value determined by Equation 1 is much higher at 0.2718 because of the smaller sample size. 552 The MANOVA performed on Donna's wordlist data produced a Pillai score of 0.1648 (p = 0.216). 553 We reiterate that sample sizes also impact *p*-values, with smaller sample sizes producing higher *p*-554 values. Taken all together, this data makes a strong case that /a/and /3/are underlyingly merged in555 Donna's speech in both speech styles.

556 On the other hand, some speakers' data are indicative of a distinction. Kim, another 48-year-557 old woman, produced 137 low back tokens in the conversational portion of her interview. Equation 558 1 suggests that a Pillai score less than 0.0397 would be evidence for a merger, but the MANOVA 559 performed on her data returned a Pillai score of 0.0658 (p = 0.010). Though the Pillai score is relatively 560 close to zero, we do not interpret her data as underlying merged, since with that many observations 561 from a truly merged distribution we would expect an even lower Pillai score (less than 0.0397). Based 562 on the 20 tokens from her wordlist, the cutoff would be 0.2718 but the MANOVA on those 20 tokens 563 yielded a Pillai score of  $0.3700 \ (p = 0.020)$ , further suggesting a distinction. Because Kim's Pillai scores 564 were higher than the thresholds and were accompanied by low *p*-values for both speech styles, we 565 conclude that Kim's low back vowels, while close in acoustic space, are not fully merged.

566 However, even when considering p-values alongside recommended thresholds, not all cases 567 are as straightforwardly interpretable as Donna's and Kim's. Scott is a 28-year-old man whose 568 interview contained 195 low back tokens. The Pillai score based on his data was 0.1975, far higher 569 than the threshold (0.0279) produced by Equation 1, and was accompanied by a low p-value (p < p570 0.001), suggesting two distinct underlying vowel distributions. However, the Pillai score based on the 571 20 tokens he produced in his wordlist (0.0397) was much lower than the threshold (again, 0.2718), 572 and had a high *p*-value (p = 0.7090). With the Pillai score lower than the threshold and accompanied 573 by a high *p*-value his wordlist data, it is tempting to interpret Scott as being a rare case of producing a 574 merger in the wordlist that he does not produce in the conversational portion of the interview. 575 However, because the sample size is so small in the wordlist, the Pillai score (and indeed, any result of 576 an inferential statistical test) should be taken with a large grain of salt. We include Scott's data here in 577 part to demonstrate that even when leveraging Pillai scores alongside a recommended threshold and 578 a *p*-value, data with low token counts can still be difficult to interpret.

579 We can add one additional tool to the suite of evidence we consider when diagnosing 580 individual speakers: a visual inspection of a plot. Figure 4 shows the distributions of  $/\alpha$  and  $/3/\beta$ 581 tokens in F1-F2 space for conversational data (left) and wordlist data (right) for both Donna (top), 582 Kim (middle), and Scott (bottom). Ellipses represent one standard deviation for each vowel class. A 583 visual inspection of these plots shows that both vowel classes exhibit a fair amount of overlap in both 584 styles, with what appears to be more overlap in Donna's, as her /3 vowel class actually encompasses 585 /a/ in both speech styles (a distributional property that indicates merger). Adding the measures of 586 Pillai scores, using the recommended thresholds for "merged" given the specific token counts, along 587 with *p*-values for these speakers in these two styles, enables us to more confidently state that Donna's 588 two vowels are underlyingly merged, while Kim's are distinct in both styles and Scott's are distinct at 589 least in the conversational style. For both Donna and Kim, the Pillai scores were higher in the wordlist 590 style compared to the conversation. Seeing these differences in Pillai scores alone, a researcher may 591 be tempted to conclude that there is style shifting occurring for both speakers, such that the merger 592 undoes itself in more careful speech. Leveraging all of the evidence together - Pillai scores alongside 593 the recommended threshold for a given sample size, as well as *p*-values and a visual inspection of the 594 data - allows us to reject this interpretation and instead see important differences between speakers 595 in the sample.







- . . . .

600 We note, in fact, that interpreting Pillai score alone without the additional evidence of 601 threshold and p-value provides a misleading interpretation of the entire dataset. Across the 30 speakers 602 in the sample who completed both tasks, a one-sided paired t-test comparing Pillai scores in 603 conversational data to Pillai scores in wordlist data suggests that there is a statistically significant trend towards higher Pillai scores (i.e., a "less merged" pronunciation) in the wordlist data (t = -3.3296, df 604 605 = 29, p = 0.001). This pattern obtains across most speakers, suggesting on the surface that almost 606 every speaker in the sample "unmerges" their vowels in wordlist style data or that they only merge the 607 vowels in casual conversation. Given that the low-back merger is a change in progress (cf. Labov, 608 Yaeger & Steiner 1972 for other examples), this pattern obtaining across this many speakers of all ages 609 is suspiciously regular – much more regular than we would expect given patterns of variation and 610 change in sociophonetic work. In fact, it was this apparently regular unmerging found in these 611 participants' wordlists that led us to investigate the effect of sample size in the first place. 612 Incorporating all of the relevant pieces of evidence (recommended threshold given sample size, along 613 with Pillai scores and p-values) allows us to understand the suspiciously regular finding as an artefact 614 of wordlists having far smaller sample sizes than conversational speech. Likewise, leveraging all of the 615 evidence, including sample size differences across styles, allows a clearer understanding of the 616 individual speakers in the sample and whether they are likely to be truly merged or not.

617

618

## K. Analysis of many speakers

619 While the level of scrutiny in the previous section is appropriate and encouraged for analyzing 620 individual speakers, we acknowledge that researchers may have a different goal in mind. For instance, 621 researchers tracking the development of a merger as it becomes closer together in phonetic space (and 622 before a categorical merging) will want to understand how the two vowel classes become less distinct 623 across generations – even before a categorical merger has taken place. Analyzing data speaker-by624 speaker requires us to take gradient measures (Pillai scores, threshold, and *p*-value) and interpret them 625 into discrete categories for each speaker and style ("merged", "not merged" or "partially merged", 626 depending on the researcher's interpretation); this discrete interpretation in turn makes it difficult to 627 track how large-scale change proceeds across generations. In this section we explore how Pillai scores 628 have changed over time in this sample by fitting a linear model to the data after implementing a 629 bootstrapping procedure to remove difference between sample sizes across speakers.

To ensure that Pillai scores are comparable across a large number of individual speakers, we recommend reducing the size of each speaker's dataset down to the sample size of the speaker who contributed the least amount of data. Downsampling in this way will allow us to obtain Pillai scores that are comparable across all speakers, and therefore allow us to track fine-grained changes in merger status across apparent time.

635 To begin our downsampling example, we first restrict our analysis to the more prolific 636 conversational portion of the interview, to maximize the possible total sample size per speaker. The least talkative speaker in this sample produced only 82 tokens in this style. So, even though some 637 speakers produced many more (as many as 327 in one case), we took a random sample (with 638 639 replacement<sup>viii</sup>) of just 82 tokens from each speaker. However, we found that Pillai scores from a 640 random sample of one speaker's data varied considerably from a different random sample from that 641 same speaker. So, we implemented a bootstrapping procedure and took 1000 random samples (with 642 replacement) of 82 tokens from each speaker's dataset, calculated the Pillai scores and other summary 643 statistics from each sample, and aggregated them by speaker. The correlation between speakers' Pillai 644 scores on the full dataset and speakers' mean Pillai scores across the 1000 samples was 0.9984, 645 suggesting that the aggregated bootstrapped values are a very good approximation of the full dataset's 646 values. The difference is that they are based on equal sample sizes rather than different sample sizes, 647 which means that the resulting Pillai scores will be comparable across speakers.



Figure 5 (color online): Log-transformed Pillai scores by gender and birth year with predicted regression lines. Lower values long the y-axis indicate a greater overlap between /a/ and /o/. The horizontal dotted line crosses the *y*-axis at log(0.0663) = -2.7136, which is based on the threshold for 82 observations calculated using Equation 1; values below that line may be considered merged.

To identify whether there were patterns across genders or time, we fit the log-transformed Pillai scores to a linear model with gender, birth year, and their interaction as predictors (Figure 5). Since the data was downsampled, Pillai scores can be directly compared and the influence of particularly talkative or reticent speakers is minimized. The results of this model show that Pillai scores for older women are low (indicating a merger) and older men are high (indicating a distinction). The difference between men and women decreases over time, with Pillai scores converging among the youngest speakers.

660 VI. CONCLUSION

648

In this paper, we present a close view into the effect of sample size on resulting Pillai scores,
a common measure for quantifying vowel overlap. We use a series of simulation experiments drawing
from an underlyingly merged pair of vowels to demonstrate (1) that larger sample sizes yield reliably

664 lower Pillai scores, (2) that unequal group sizes across the two vowel classes is irrelevant in the 665 calculation of Pillai scores, and (3) that it takes more data than many sociolinguistic studies collect to 666 reliably return a low Pillai score (e.g., under 0.1) even for underlyingly merged data. These results have 667 implications for how Pillai scores are compared across studies and between speakers or speech styles 668 within the same study. We provide some recommendations for maximizing reliability in the use of 669 Pillai scores, and provide a formula to assist researchers in determining a reasonable Pillai score 670 threshold to use as an indicator of merged status given their sample size. We recommend the use of Equation 1 a in conjunction with the Pillai scores' accompanying *p*-values to make informed decisions 671 about the merged status of two vowels in a given speaker. By properly using and reporting aspects of 672 673 Pillai score, researchers can come closer to accurately quantifying vowel overlap, identifying vowel 674 mergers, and ultimately understanding broad patterns of variation and change in vowel merger.

675

676 ACI

### ACKNOWLEDGMENTS

Thanks to the Linguistics Discussion Group at BYU and the audience at ASA2021 in Seattle for their useful comments. We thank William Christensen and Joe Fruehwald for their help in understanding underlying distributions and transformations. We also thank the two anonymous reviewers for their thoughtful and encouraging feedback, without which this paper would have been substantially less clear and less impactful. Any remaining faults are our own. The first author also graciously acknowledges the University of Georgia Graduate School Dean's Award for funding the fieldwork that produced the data used in the case study.

- 686 See supplementary material at [URL will be inserted by AIP] for a brief tutorial on how to
- 687 implement these recommendations in R.
- 688
- 689
- 690 **REFERENCES**
- Amengual, Mark & Pilar Chamorro. 2015. The Effects of Language Dominance in the Perception
   and Production of the Galician Mid Vowel Contrasts. *Phonetica* 72(4). 207–236.
   https://doi.org/10.1159/000439406.
- Arnold, Jeffrey B. 2018. ggthemes: Extra Themes, Scales and Geoms for "ggplot2."
   https://CRAN.R-project.org/package=ggthemes.
- Austin, Martha. 2020. Mismatches between Linguistic and Sociolinguistic Perception. Presented
   at the The 94th Annual Meeting of the Linguistic Society of America, New Orleans, LA.
- Bartlett, M. S. 1939. A note on tests of significance in multivariate analysis. *Mathematical Proceedings of the Cambridge Philosophical Society* 35(2). 180–185.
   https://doi.org/10.1017/S0305004100020880.
- Becker, Kara. 2019. Introduction. In Kara Becker (ed.), *The Low-Back-Merger Shift: Uniting the Canadian Vowel Shift, the California Vowel Shift, and short front vowel shifts across North America* (Publication of the American Dialect Society 104). Durham, NC: Duke
   University Press.
- Berry, Grant M & Mirjam Ernestus. 2018. Phonetic alignment in English as a *lingua franca*:
   Coming together while splitting apart. *Second Language Research* 34(3). 343–370.
   https://doi.org/10.1177/0267658317737348.
- Berry, Grant Michael. 2018. *Liminal voices, central constraints: Minority adoption of majority sound change*. The Pennsylvania State University Ph.D. Dissertation.
- Bray, Andrew. 2021. ['hqki]: An Emerging Third-Order Index of a Hockey-Based Persona.
   Presented at the New Ways of Analyzing Variation 49, Austin, Texas.
- Bray, James H. & Scott E. Maxwell. 1985. *Multivariate Analysis of Variance* (Quantitative
   Applications in the Social Sciences). Vol. 07–054. Newbury Park, CA: Sage Publications.
- Brozovsky, Erica Sharon. 2020. *Taiwanese Texans: A Sociolinguistic Study of Language and Cultural Identity*. Austin, TX: University of Texas at Austin Ph.D. Dissertation.
- Di Paolo, Marianna & Alice Faber. 1990. Phonation differences and the phonetic content of the
   tense-lax contrast in Utah English. *Language Variation and Change* 2(02). 155–204.
   https://doi.org/10.1017/S0954394500000326.
- Di Paolo, Marianna, Malcah Yaeger-Dror & Alicia Beckford Wassink. 2011. Analyzing Vowels. In
   Marianna Di Paolo & Malcah Yaeger-Dror (eds.), *Sociophonetics: A Student's Guide*, 87–
   106. 1st edn. London: Routledge.
- Fieberg, John & Christopher O. Kochanny. 2005. Quantifying Home-Range Overlap: The
   Importance of the Utilization Distribution. (Ed.) Lanham. *Journal of Wildlife*

724	Management 69(4). 1346–1359. https://doi.org/10.2193/0022-
725	541X(2005)69[1346:QHOTIO]2.0.CO;2.
726	Fisher, Sabriya, Hilary Prichard & Betsy Sneller. 2015. The Apple Doesn't Fall Far From the Tree:
727	Incremental Change in Philadelphia Families. University of Pennsylvania Working Papers
728	in Linguistics 21(2). http://repository.upenn.edu/pwpl/vol21/iss2/7.
729	Freeman, Valerie. 2021. Vague eggs and tags: Prevelar merger in Seattle. Language Variation
730	and change 1–24. https://doi.org/doi:10.1017/S0954394521000028.
731	Fung, Roxana S. Y. & Chris K. C. Lee. 2019. Tone mergers in Hong Kong Cantonese: An
732	asymmetry of production and perception. The Journal of the Acoustical Society of
733	America 146(5). EL424–EL430. https://doi.org/10.1121/1.5133661.
734	Gorman, Kyle & Daniel Ezra Johnson. 2013. Quantitative Analysis. In Robert Bayley, Richard
735	Cameron & Ceil Lucas (eds.), The Oxford Handbook of Sociolinguistics, vol. 1. Oxford
736	University Press. https://doi.org/10.1093/oxfordhb/9780199744084.013.0011.
737	Hall-Lew, Lauren. 2013. 'Flip-flop' and mergers-in-progress. English Language and Linguistics
738	17(02). 359–390.
739	Hall-Lew, Lauren, Ruth Friskney & James M. Scobbie. 2017. Accommodation or political identity:
740	Scottish members of the UK Parliament. Language Variation and Change 29(3). 341–
741	363. https://doi.org/10.1017/S0954394517000175.
742	Han, Jeong-Im & Hyunsook Kang. 2013. Cross-generational Change of /o/ and /u/ in Seoul
743	Korean I: Proximity in Vowel Space. Phonetics and Speech Sciences 5(2). 25–31.
744	https://doi.org/10.13064/KSSS.2013.5.2.025.
745	Havenhill, Jonathan. 2015. Maintenance of the COT-CAUGHT Contrast Among Metro Detroit
746	Speakers: A Multimodal Articulatory Analysis. University of Pennsylvania Working Papers
747	in Linguistics 21(2).
748	Hay, Jennifer, Paul Warren & Katie Drager. 2006. Factors influencing speech perception in the
749	context of a merger-in-progress. Journal of Phonetics (Modelling Sociophonetic
750	Variation) 34(4). 458–484. https://doi.org/10.1016/j.wocn.2005.10.001.
751	Holland, Cory & Tara Brandenburg. 2017. Beyond the Front Range: The Coloradan Vowel Space.
752	In Valerie Fridland, Alicia Beckford Wassink, Tyler Kendall & Besty E. Evans (eds.), Speech
753	in the Western States, Volume 2: The Mountain West (Publication of the American
754	Dialect Society 102), 9–30. Durham, NC: Duke University Press. DOI: 10.1215/00031283-
755	4295277.
756	Islam, Md Jahurul & Iftakhar Ahmed. 2020. Mid-front and back vowel mergers in Mymensingh
757	Bangla: An acoustic investigation. <i>Linguistics Journal</i> 14(1). 206–232.
758	Ismay, Chester & Albert Y. Kim. 2020. Statistical Inference via Data Science: A ModernDive into
759	R and the Tidyverse (The R Series). Boca Raton, FL: Taylor & Francis Group.
760	https://moderndive.com/index.html.
761	Jensen, Sandra & Natalie Braber. 2021. The BATH-TRAP split in the East Midlands. Poster
762	presented at the New Ways of Analyzing Variation 49, Austin, Texas.
763	Jibson, Jonathan. 2021. Merged status thresholds for Pillai scores. Poster presented at the New
764	Ways of Analyzing Variation 49, Austin, Texas.
765	Johnson, Daniel Ezra. 2010. Stability and Change Along a Dialect Boundary: The Low Vowels of
766	Southeastern New England (Publication of the American Dialect Society 95). Durham,
767	NC: Duke University Press.

- Johnson, Daniel Ezra. 2015. Quantifying vowel overlap with Bhattacharyya's affinity. Presented
   at the New Ways of Analyzing Variation (NWAV44), Toronto.
- Johnson, Richard A. & Dean W. Wichern. 2012. *Applied Multivariate Statistical Analysis*. Phi
   Learning Private Limited.
- Joo, Hyoun-A, Lara Schwarz & B. Richard Page. 2018. Nonconvergence and Divergence in
  Bilingual Phonological and Phonetic Systems: Low Back Vowels in Moundridge
  Schweitzer German and English. *Journal of Language Contact* 11(2). 304–323.
  https://doi.org/10.1163/19552629-01102006.
- Kelley, Matthew C. & Benjamin V. Tucker. 2020. A comparison of four vowel overlap measures.
   *The Journal of the Acoustical Society of America* 147(1). 137–145.
   https://doi.org/10.1121/10.0000494.
- Kendall, Tyler & Valerie Fridland. 2017. Regional relationships among the low vowels of U.S.
   English: Evidence from production and perception. *Language Variation and Change* 29(2). 245–271. https://doi.org/10.1017/S0954394517000084.
- Kendall, Tyler & Valerie Fridland. 2021. *Sociophonetics* (Key Topics in Sociolinguistics).
   Cambridge: Cambridge University Press.
- 784 Kettig, Thomas T. 2021. *Ha'ina 'ia Mai Ana Ka Puana: The Vowels of 'Ōlelo Hawai'i*. Mānoa,
  785 Hawai'i: University of Hawai'i at Mānoa Ph.D. Dissertation.
- Labov, William. 1994. Principles of linguistic change. Vol. 1: Internal features (Language in
   Society). Oxford: Wiley-Blackwell.
- Labov, William & Maciej Baranowski. 2006. 50 msec. Language Variation and Change 18(03).
   https://doi.org/10.1017/S095439450606011X.
- Labov, William, Sabriya Fisher, Duna Gylfadottír, Anita Henderson & Betsy Sneller. 2016.
  Competing systems in Philadelphia phonology. *Language Variation and Change* 28(3).
  273–305. https://doi.org/10.1017/S0954394516000132.
- Labov, William, Malcah Yaeger & Richard Steiner. 1972. A quantitative study of sound change in progress: Volume 1. Philadelphia, PA: US Regional Survey.
- Lüdecke, Daniel, Mattan Ben-Shachar, Indrajeet Patil, Philip Waggoner & Dominique Makowski.
   2021. performance: An R Package for Assessment, Comparison and Testing of Statistical
   Models. *Journal of Open Source Software* 6(60). 3139.
- 798 https://doi.org/10.21105/joss.03139.
- Lüdecke, Daniel, Indrajeet Patil, Mattan Ben-Shachar, Brenton Wiernik, Philip Waggoner &
   Dominique Makowski. 2021. see: An R Package for Visualizing Statistical Models. *Journal of Open Source Software* 6(64). 3393. https://doi.org/10.21105/joss.03393.
- Metropolis, Nicholas & S. Ulam. 1949. The Monte Carlo Method. *Journal of the American Statistical Association* 44(247). 335–341.
- 804 https://doi.org/10.1080/01621459.1949.10483310.
- Moulton, William G. 1968. Structural Dialectology. *Language* 44(3). 17.
- Nadeu, Marianna & Margaret E.L. Renwick. 2016. Variation in the lexical distribution and
   implementation of phonetically similar phonemes in Catalan. *Journal of Phonetics* 58.
   22–47. https://doi.org/10.1016/j.wocn.2016.05.003.
- Newbert, Cornelia. 2021. Language variation in South Africa: A sociophonetic study of the vowel
   system of Black South African English. Regensburg, Germany: University of Regensburg
   Ph.D. Dissertation.

812 Nycz, Jennifer & Lauren Hall-Lew. 2013. Best practices in measuring vowel merger. Proceedings 813 of Meetings on Acoustics 20(1). 060008. https://doi.org/10.1121/1.4894063. 814 Olson, Chester L. 1976. On choosing a test statistic in multivariate analysis of variance. 815 Psychological Bulletin 83(4). 579–586. https://doi.org/10.1037/0033-2909.83.4.579. 816 Pedersen, Thomas Lin & Fabio Crameri. 2020. scico: Colour Palettes Based on the Scientific Colour-Maps. https://CRAN.R-project.org/package=scico. 817 818 Pfiffner, Alexandra M. 2021. Cue-Based Features: Modeling change and varitaion in the voicing 819 contrasts of Minnesotan English, Afrikaans, and Dutch. Washington, DC: Georgetown 820 University Ph.D. Dissertation. Pillai, K. C. S. 1955. Some New Test Criteria in Multivariate Analysis. The Annals of Mathematical 821 822 Statistics 26(1). 117–121. https://doi.org/doi:10.1214/aoms/1177728599. 823 R Core Team. 2021. R: A Language and Environment for Statistical Computing. Vienna, Austria: 824 R Foundation for Statistical Computing. http://www.R-project.org. 825 Regan, Brendan. 2020. Extending Pillai Scores to Fricative Mergers: Advancing a Gradient 826 Analysis of a Split-in-Progress in Andalusian Spanish. University of Pennsylvania Working 827 Papers in Linguistics 26(2). Rencher, Alvin C. & William F. Christensen. 2012. *Methods of multivariate analysis* (Wiley Series 828 829 in Probability and Statistics). Third Edition. Hoboken, New Jersey: Wiley. 830 Schmidt, Penelope, Chloé Diskin-Holdaway & Debbie Loakes. 2021. New insights into /el/-/æl/ 831 merging in Australian English. Australian Journal of Linguistics 41(1). 66–95. 832 https://doi.org/10.1080/07268602.2021.1905607. 833 Seaman, D. Erran, Joshua J. Millspaugh, Brian J. Kernohan, Gary C. Brundige, Kenneth J. 834 Raedeke & Robert A. Gitzen. 1999. Effects of Sample Size on Kernel Home Range 835 Estimates. The Journal of Wildlife Management 63(2). 739. 836 https://doi.org/10.2307/3802664. 837 Sloos, Marjoleine. 2013. The reversal of the BÄREN-BEEREN merger in Austrian Standard 838 German. The Mental Lexicon 8(3). 353–371. https://doi.org/10.1075/ml.8.3.05slo. Sneller, Betsy. 2018. Mechanisms Of Phonological Change. Philadelphia: University of 839 Pennsylvania Ph.D. Dissertation. 840 Stanford, James N., Monica Nesbitt, James King & Sebastian Turner. 2019. Pioneering a dialect 841 842 shift in the Pioneer Valley: Evidence for the Low-Back-Merger Shift in Western 843 Massachusetts. Presented at the New Ways of Analyzing Variation 48, Eugene, Oregon. 844 Stanley, Joseph A. 2020. Vowel dynamics of the Elsewhere Shift: A sociophonetic analysis of 845 English in Cowlitz County, Washington. Athens, Georgia: University of Georgia Ph.D. 846 Dissertation. 847 Stanley, Joseph A. 2021. joeyr: Functions for Vowel Data. https://joeystanley.github.io/joeyr/. 848 Strelluf, Christopher. 2018. Speaking from the heartland: The midland vowel system of Kansas 849 City (Publication of the American Dialect Society 103). Durham, NC: Duke University 850 Press. 851 Tse, Holman. 2018. Beyond the Monolingual Core and out into the Wild: A Variationist Study of 852 Early Bilingualism and Sound Change in Toronto Heritage Cantonese. Pittsburgh, PA: 853 University of Pittsburgh Ph.D. Dissertation. 854 Upton, Graham & Ian Cook. 2014. A Dictionary of Statistics. Oxford University Press. 855 https://doi.org/10.1093/acref/9780199679188.001.0001.

- Venables, W. N., Brian D. Ripley & W. N. Venables. 2002. *Modern applied statistics with S*(Statistics and Computing). 4th ed. New York: Springer.
- Warner, Natasha, Allard Jongman, Joan Sereno & Rachèl Kemps. 2004. Incomplete
   neutralization and other sub-phonemic durational differences in production and
   perception: evidence from Dutch. *Journal of Phonetics* 32(2). 251–276.
- 861 https://doi.org/10.1016/S0095-4470(03)00032-9.
- Warren, Paul. 2018. Quality and quantity in New Zealand English vowel contrasts. *Journal of the International Phonetic Association* 48(3). 305–330.
- 864 https://doi.org/10.1017/S0025100317000329.
- Wasserstein, Ronald L., Allen L. Schirm & Nicole A. Lazar. 2019. Moving to a World Beyond "p <</li>
  0.05." *The American Statistician* 73(sup1). 1–19.
- 867 https://doi.org/10.1080/00031305.2019.1583913.
- Wassink, Alicia Beckford. 2006. A geometric representation of spectral and temporal vowel
   features: Quantification of vowel overlap in three linguistic varieties. *The Journal of the Acoustical Society of America* 119(4). 2334–2350.
- Whalen, D. H. & Wei-Rong Chen. 2019. Variability and Central Tendencies in Speech Production.
   *Frontiers in Communication* 4. 49. https://doi.org/10.3389/fcomm.2019.00049.
- Wickham, Hadley. 2015. ggplot2: Elegant Graphics for Data Analysis (Use R!). 2nd edn. New
  York: Springer.
- Wickham, Hadley, Mara Averick, Jennifer Bryan, Winston Chang, Lucy McGowan, Romain
   François, Garrett Grolemund, et al. 2019. Welcome to the Tidyverse. *Journal of Open Source Software* 4(43). 1686. https://doi.org/10.21105/joss.01686.
- Wong, Amy Wing-mei & Lauren Hall-Lew. 2014. Regional variability and ethnic identity: Chinese
  Americans in New York City and San Francisco. *Language & Communication* 35. 27–42.
  https://doi.org/10.1016/j.langcom.2013.11.003.
- 881

<sup>i</sup> See also Kelley & Tucker (2020) for details on the overlapping coefficient, another nonparametric method for calculating overlap.

<sup>ii</sup> Preliminary work for Sneller (2018) attempted, with mixed results, a modified model approach of Bhattaharyya's Affinity by extracting model values for fixed effects and manually adjusting the data. However, the onerousness of this approach makes it not widely implementable in comparison with Pillai, which has the advantage of being able to easily integrate commonly used mixed-effects models.

<sup>iii</sup> For the uncorrelated data, the bivariate normal distribution was most easily generated by combining two independent univariate normal distributions, because the product of their probability densities is equal to their joint probability densities (Johnson & Wichern 2012 chapter 4 page 151). We generated F1 in R by running rnorm(x, mean =0, sd=1), where x is the number of tokens generated. F2 was generated using the same code, and the two sets of numbers were combined to create the bivariate normal distribution.

<sup>iv</sup> In this paper, we focus on the effect of sample size for underlyingly merged speakers. Future work may fruitfully investigate how sample size interacts with Pillai score for underlyingly unmerged speakers as well.

<sup>v</sup> Let  $p_{95} = e^{1-\log(m)}$ . Since  $e^{x-y} = e^{x+(-y)} = e^x e^{-y} = e^x \frac{1}{e^y} = \frac{e^x}{e^y}$ , and since  $e^1 = e$ , then  $p_{95} = \frac{e}{e^{\ln(m)}}$ . Since  $e^{\ln(x)} = x$ , then  $p_{95} = \frac{e}{m}$ . (We are grateful to the anonymous review who pointed out these simplification steps for this formula.) This is most easily implemented in R as exp(1)/m.

<sup>vi</sup> We also reiterate the caution about small sample sizes overall, and point out that very small samples may increase the likelihood of a Type II error (i.e, a false negative or failing to reject the null hypothesis when it is actually false), leading the researcher to conclude that two vowels are merged when in reality there is just not enough data to detect a distinction.

<sup>vii</sup> One reviewer asked why we did not recommend accounting for the difference in sample sizes across speech styles by simply adding a term for style in the linear model. The reasoning is fairly

straightforward: researchers cannot control how many tokens of interest are produced in the conversational portion of the interview, allowing the difference in token count across style to vary wildly by speaker. Since sample sizes influence Pillai score in a predictable way but style influences Pillai score in an unpredictable way (as it's dependent on the differences in sample sizes), it is inadvisable to use style as the predictor rather than a more straightforward way to control for or incorporate sample size.

viii Note that we sample with replacement (replace=TRUE), in order to make our resampling comparable across all speakers including our least talkative speaker.

If we were to resample *without* replacement, this introduces a confound related to sample size: the amount of error introduced per speaker is proportional to their sample size. Resampling without replacement 1000 times would produce 1000 identical distributions for the least talkative speaker, but 1000 different distributions for every other speaker. Resampling *with* replacement allows us to obtain a standard deviation for each speaker with similar confidence, and introduces a similar amount of uncertainty across speakers in the means of their distributions. For more information on bootstrapping with replacement, we refer the reader to Ismay & Kim (2020).