

# Discriminating speakers using perceptual clustering interface

Benjamin O'Brien<sup>1</sup>, Alain Ghio<sup>1</sup>, Corinne Fredouille<sup>2</sup>, Jean-François Bonastre<sup>2</sup>, Christine Meunier<sup>1</sup>

<sup>1</sup> Aix-Marseille Univ., CNRS, LPL, UMR 7309, Aix-en-Provence, FR

<sup>2</sup> Laboratoire d'Informatique d'Avignon, Avignon Université, Avignon, FR

**INTRODUCTION** The challenges facing listeners tasked to identify speakers are well documented.<sup>1 2 3</sup> In addition to providing listeners with high-quality speech recordings that accurately represent the speakers, the method of presentation itself is equally important.<sup>4 5</sup> Numerous perception studies have employed a binary approach, where participants are asked to judge whether two speech recordings are similar or different, as a way of examining the effects of such things as noise,<sup>6</sup> language familiarity,<sup>7 8</sup> and stimuli selection methods.<sup>9</sup> Oftentimes this requires numerous tests, which can be time-consuming for participants and experimenters. Moreover, there persists concern for memory bias, as a “fresh” voice is not equivalent to a voice that was presented in a previous binary test.

As an alternative, we proposed the development of a perceptual *clustering* method, which is often employed in the domain of machine learning.<sup>10 11</sup> We theorized that this approach would allow users to better personalize their engagements with speech materials and organize their proximities in relation to their perceived likeness. In addition, it was more economical in terms of the number of trials required to assess a listener’s ability to identify speakers.

In order to study the speaker discrimination performance of participants using a perceptual clustering interface, it was important to organize and select stimuli based on how listeners perceive them as similar or different. Studies suggest listeners rely on a common set of acoustic features to identify speakers.<sup>12 13</sup> It is common in the development of automatic voice recognition and speaker identification system to extract MFCCs from speech recordings to train models. A popular trend in the field involves the transformation of these features into i-vectors, which have been shown to be quite accurate in identifying speakers.<sup>14 15</sup> Recent work has shown that Cosine Distance Scoring (CDS) with Within-Class covariance normalization (WCCM) is similarly effective while reducing the complexity of the task.<sup>16</sup> Our second objective was to examine the relationships between participant performance and the CDS generated from the speaker i-vectors.

**METHODS** Speech recordings were selected from the PTSTVox database,<sup>17</sup> which included 24 francophone speakers (12 female, 12 male) who recited three French-texts into a Zoom H4N stereo microphone (sampling rate: 44.1 kHz; bit depth: 16-bit) over the course of two recording sessions (mean  $118.96 \pm 17.54$  s). SPro<sup>18</sup> was used to extract 19 MFCCs, deltas, and delta-deltas from each recording. ALIZE<sup>19</sup> was used to compress these features into i-vectors and then calculate CDS between each one, whereupon the WCCM was computed over the entire set. Two groups of five speakers were selected: the *Alpha* group was composed of speakers with the greatest distance between them and the *Betha* group was composed of speakers with the smallest distance between them. For each speaker, twelve utterances were selected (120 recordings; mean  $1.47 \pm 0.51$  s). Groups were divided into six sessions, such that each session was balanced and composed of four different (non-repeating) chunks per speaker.

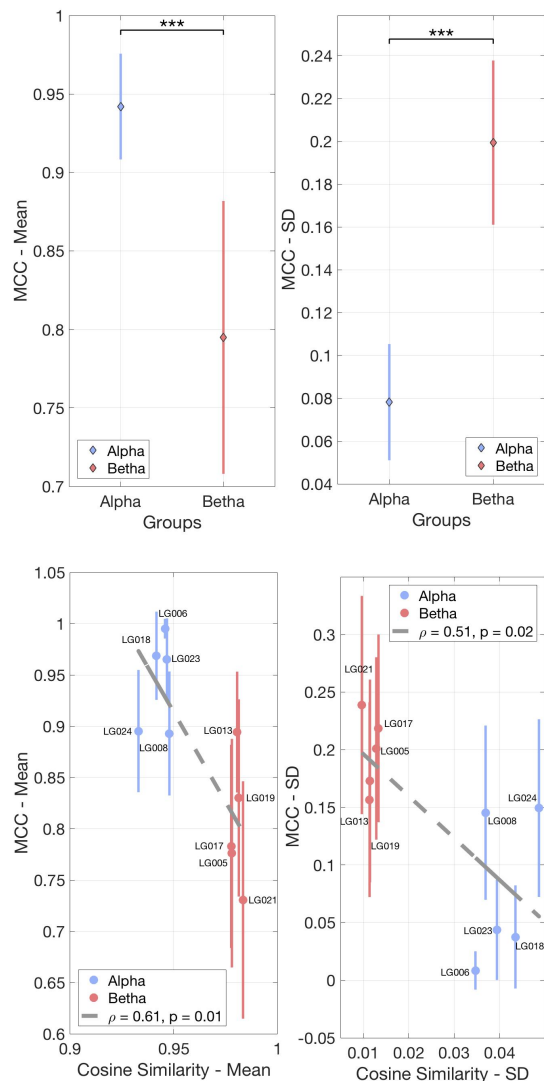
Twenty-four people (14 female;  $24.2 \pm 6.7$  years), who self-reported good hearing, participated in our study. Their task was to group 20 speech recordings into five cluster groups, where each cluster represented a unique speaker. To do this, they used the TCL-LABX interface,<sup>20</sup> which allowed them to move recordings in a 2-D space and assign them to different clusters. They completed six sessions.

The Mathews Correlation Coefficient (MCC) was selected to determine how accurate the participants were at discriminating speakers (1), where *TP*, *TN*, *FP*, *FN* represent the selections that were “true positive,” “true negative,” “false positive,” and “false negative,” respectively. The mode speaker in each cluster was used to calculate the MCC. MCC means and standard deviations were calculated for each speaker.

$$MCC = \frac{(TP*TN)-(FP*FN)}{\sqrt{((TP+FP)*(TP+FN))*(TN+FP)*(TN+FN)}} \quad (1)$$

**RESULTS** To examine participant performance discriminating speakers, two-level nested ANOVA procedures were applied to MCC mean and standard deviation for groups with different speakers. We found a main effect on groups for MCC mean  $F_{1,240} = 32.92$ ,  $p < 0.001$ ,  $\eta_p^2 = 0.12$ , and no significant differences between speakers within each group,  $p > 0.05$ . Post-hoc tests revealed Alpha had a higher MCC mean ( $0.94 \pm 0.20$ ) when compared to Betha ( $0.8 \pm 0.02$ ),  $p < 0.001$ . Similarly we found a main effect on MCC standard deviation  $F_{1,240} = 26.04$ ,  $p < 0.001$ ,  $\eta_p^2 = 0.1$ , but again no significant differences between speakers

within each group,  $p > 0.05$ . Post-hoc tests revealed Alpha had a lower MCC standard deviation ( $0.08 \pm 0.02$ ) when compared to Beta ( $0.2 \pm 0.02$ ),  $p < 0.001$  (Fig. 1a).



**Figure 1a:** Mean (Left) and standard deviation (Right) of participant MCC per group. Diamonds and vertical lines represent the means and standard errors, respectively. {\*\*\*} signifies  $p < 0.001$ . **Figure 1b:** Pearson's coefficient was used to examine the relationship between CDS and MCC metrics: mean (Left)  $\rho = 0.61, p = 0.01$ , and standard deviation (Right)  $\rho = 0.51, p = 0.02$ . Circles and vertical lines represent the means and standard errors, respectively. The text indicates the speaker id.

We then examined whether our method of selecting and grouping speakers played a role in participant performance. For each speaker we calculated the CDS mean and standard deviation between it and the other group speakers and then calculated the Pearson's correlation coefficient to examine the relationships between the two metrics. The speaker CDS mean difference estimated the MCC mean at  $\rho = 0.61, p = 0.01$ , whereas the speaker CDS standard deviation estimated the MCC standard deviation at  $\rho = 0.51, p = 0.02$  (Fig. 1-b).

**DISCUSSION** This study demonstrated that users were able to use a clustering interface to make discriminations based on their perceived differences between speech recordings. Participants performed at a relatively high level, as indicated by the mean and standard MCC values, which suggests they found the interface easy to navigate and efficient to use. In addition, the significant differences between groups also underscore the importance of developing methods for selecting and grouping speakers. We observed that as the CDS mean increased, participants were less accurate discriminating speakers, and, conversely, as the CDS standard deviation decreased, participants showed greater variability. These findings have led us to develop a new study to compare the effects of presentation on users performing speaker discrimination tasks with similar speaker stimuli.

## References

- Cambier-Langeveld, T., Rossum, M., Vermeulen, J. (2014). Whose voice is that? Challenges in forensic phonetics.
- Mattys, S. Davis, M., Bradlow, A., Scott, S. (2012). Speech recognition in adverse conditions: A review. *Language and Cognitive Processes - LANG COGNITIVE PROCESS*. 27. 953-978.
- Nolan, F. (2001). "Speaker identification evidence: its forms, limitations, and roles." *Law and Language: Prospect and Retrospect*.
- Boë, L., Bonastre, J-F. (2012). L'identification du locuteur: 20 ans de témoignage dans les cours de Justice. Le cas du LIPSADON « laboratoire indépendant de police scientifique », JEP, Grenoble: 417-424.
- Hollien, H., Bahr, R., Künzel, H., Hollien, P. (2013). "Criteria for earwitness lineups," *Int. Jnl of Speech Language and the Law* 2, 143-153.
- Smith, H., Baguley, T., Robson, J., Dunn, A., Stacey, P. (2018). Forensic voice discrimination: The effect of speech type and background noise on performance. *Applied Cognitive Psychology*.
- Fleming, D., Giordano, B. L., Caldara, R., & Belin, P. (2014). A language familiarity effect for speaker discrimination without comprehension. *Proc National Academy of Sciences*, 111(38)
- Levi, S. V., & Schwartz, R. G. (2013). The development of language specific and language-independent talker processing. *Journal of Speech, Language, and Hearing Research*, 56(3), 913-920.
- Mühl, C., Sheil, O., Jarutytė, L. et al. (2018) The Bangor Voice Matching Test: A standardized test for the assessment of voice perception ability. *Behav Res* 50, 2184-2192.
- Kinnunen, T. and Kilpeläinen, T. (2000). Comparison of clustering algorithms in speaker identification. *Proc. IASTED Int. Conf. Signal Processing and Communications (SPC)*. 222-227.
- Lukic, Y., Vogt, C., Dürr, O., Stadelmann, T. (2016). Speaker identification and clustering using convolutional neural networks. 1-6.
- LaRivière, C. (1971). "Some acoustic and perceptual correlates of speaker identification," *Proc 7<sup>th</sup> Int. Congress Phonetic Sciences*: 558-564.
- Roebuck, R., and Wilding, J. (1993). "Effects of vowel variety and sample length on identification of a speaker in a line-up," *Applied Cognitive Psychology* 7: 475-481.
- Dehak, N., Kenny, P., Dehak, R., Dumouchel, P., Ouellet, P. (2011) Front-End Factor Analysis for Speaker Verification, in *IEEE Transactions on Audio, Speech, and Language Processing*, 19(4): 788-798.
- Kanagasundaram, A., Vogt, R., Dean, D., Sridharan, S., Mason, M. (2011). i-vector Based Speaker Recognition on Short Utterances. *Proc International Speech Communication Association, INTERSPEECH*.
- Fredouille, C., Charlet, D. (2014) Analysis of I-Vector framework for Speaker Identification in TV-shows. *Interspeech*, Singapore, Singapore.
- Chanclu, A., Georgeton, L., Fredouille, C., Bonastre, J-F. (2020) PTSVOX: une base de données pour la comparaison de voix dans le cadre judiciaire. 6e conférence conjointe Journées d'Études sur la Parole, Nancy, FR. pp.73-81
- Speech Signal Processing (SPro) Toolkit. <https://www.irisa.fr/metiss/guig/spro>
- Larcher, A., Bonastre, J-F., Fauve, B, et al. (2013) "ALIZE 3.0 - Open source toolkit for state-of-the-art speaker recognition." *Interspeech*, Lyon.
- Gaillard, P. (2009). Laissez-nous trier ! TCL-LabX et les tâches de catégorisation libre de sons.