



Seeking voice twins – an exploration of VoxCeleb using automatic speaker recognition and two clustering methods

Linda Gerlach^{1,2} (Ig589@cam.ac.uk), Kirsty McDougall¹, Finnian Kelly², and Anil Alexander² ¹University of Cambridge, Cambridge, UK ²Oxford Wave Research Ltd., Oxford, U.K.

1. Motivation

What are voice twins?

- Located at the extreme end of perceived voice similarity, voice twins may be different, unrelated speakers that sound extremely similar to one another.
- Acoustic features contributing to speaker similarity include F0, LTF1 to LTF4 [1].

Research questions:

Is it possible to find different, unrelated speakers that sound extremely similar to another, i.e. voice twins, one

2. Proposed methodology

Finding speakers that sound extremely similar:

- I. Select perceptually relevant features
- 2. Model features automatically
- 3. Calculate similarity estimates for speaker pairs

Why should we care?

• Voice twins could help in gaining further insights into human perception of speaker individuality.





- automatically?
- Can listeners successfully human differentiate between speech from "voice twins", from the same speakers, and from random speaker pairs?
- How confident are listeners in their judgements?
- 4. Use these similarity estimates to find subgroups of similar-sounding speakers
- 5. Explore clustering techniques to find these subgroups of potentially indistinguishable speakers
- 6. Assess how difficult it is to differentiate between speakers in the subgroups identified using human listeners

3. Speaker databases

- Subsets from three speaker databases: small and controlled to large and diverse
- Selection of good quality recordings:
 - Net speech > 20 s
 - \odot SNR > 18 dB
- $\star \star \star \star \star \star$
- Clipping < 50 %
- Selection of two random files per speaker to model intra- and interspeaker variability

WYRED [2]: Homogeneous database 180 male speakers of similar age

4. Experiment

Automatic speaker recognition

• Extraction of perceptually relevant phonetic features LTF1 to LTF4 using x-vectors [5]



• Creation of voice models and calculation of similarity scores using VOCALISE automatic speaker recognition software [6]



English language, similar accents Task 2 (accomplice telephone call) Task 4 (answerphone message) Studio quality recordings

GBR-ENG [3]: 390 speakers (202 female) English language, various accents Landline recordings made in GB

VoxCeleb [4]:



Highly diverse database 906 speakers (metadata not available) Mostly English language, various accents Interview recordings from YouTube

Agglomerative hierarchical clustering (AHC) groups recordings hierarchically according to pairwise distances (similarity) [7].

- Advantage:
 - All recordings will form part of a cluster
- Disadvantages:
 - Clusters are determined visually in dendrogram
 - Maximum number of clusters needed as input
- Pilot experiment using a GBRENG subset [8]:
 - 180 unique speaker recordings (90 female)
 - Clusters according to speaker sex
 - Difficult to assess cluster similarity further

DBSCAN (Density-Based Spatial Clustering of Applications with Noise) groups recordings within dense areas as defined by a minimum number of neighbouring points within the radius EPS [9].

- Advantages:
 - Recordings that are not close to others are penalised (= noise)
 - Varying the distance measure EPS may result in voice pairs of different degrees of similarity
- Disadvantage:
 - Not all recordings may form part of a cluster
- In this experiment:
 - Male and female speakers processed together
 - Voice twin candidates = clusters containing both recordings of two different speakers

5. Current results

Number of voice twin candidates per database for different distance values (EPS):

6. Next: Listener experiment

- Three blocks (one per dataset)
- Same/different judgements, 1-5 confidence scale
- Sample duration: 3-5s
- Comparisons to include:
 - Voice twins



7. Hypotheses assessed

- The automatic method produces voice twin pairs that can be used in listener experiments.
- perform Listeners will when worse 2. discriminating between speech from voice twins than from random speakers (higher

EPS	WIKED	GDREING	voxCeleb
0.2	1	3 ♀, 1 ♂	1 ♀,0♂
0.25	4	4 9,6 🔿	19,10
0.3	4	2 🖓 , 8 🔿	0 ♀, 2 ♂

- Same-speaker pairs
- Random different-speaker pairs
- Two files per speaker, three samples per recording
- Listeners will judge subsets of the total number of comparisons

Ask to listen to some samples! number of false acceptances).

Listeners' confidence in their judgements will vary according to the type of comparison in line with [10].

References

[1] McDougall, K. (2021). Ear-Catching versus Eye-Catching? Some Developments and Current Challenges in Earwitness Identification Evidence. Proceedings of XVII AISV (Associazione Italiana Scienze Della Voce) Conference: 'Speaker Individuality in Phonetics and Speech Sciences: Speech Technology and Forensic Applications', [Preprint]. [2] Gold, E., Ross, S., & Earnshaw, K. (2018). The 'West Yorkshire Regional English Database': Investigations into the generalizability of reference populations for forensic speaker comparison casework. Proceedings of the Annual Conference of the International Speech Communication Association, INTERSPEECH, Paper 0065, 2748-2752. [3] GBR-ENG database. (2019). A telephonic speech database collected for the UK Government for evaluating speech technologies. [4] Nagrani, A., Chung, J. S., Xie, W., & Zisserman, A. (2020). Voxceleb: Large-scale speaker verification in the wild. Computer Speech & Language, 60, 101027. [5] Gerlach, L., McDougall, K., Kelly, F., & Alexander, A. (2021, August). How do automatic speaker recognition systems 'perceive' voice similarity? Further exploration of the relationship between human and machine voice similarity ratings. Proceedings of the Annual Conference of the International Association for Forensic Phonetics and Acoustics (IAFPA).

[6] Kelly, F., Forth, O., Kent, S., Gerlach, L., Alexander, A. (2019). Deep neural network based forensic automatic speaker recognition in VOCALISE using x-vectors. Audio Engineering Society (AES) Forensics Conference 2019, Porto, Portugal.

[7] San Segundo, E., Foulkes, P., French, P., Harrison, P., Hughes, V., & Kavanagh, C. (2018). Cluster analysis of voice quality ratings: Identifying groups of perceptually similar speakers. In: Proceedings of the Conference on Phonetics & Phonology in German-speaking countries (P&P 13). pp.173-176. [8] Gerlach, L., McDougall, K., Kelly, F., & Alexander, A. (2022, April). Selecting similar-sounding speakers for forensic phonetic applications: an exploration of cluster analysis using automatic speaker recognition. Presented at British Association of Academic Phoneticians (BAAP) Colloquium. [9] Schubert, E., Sander, J., Ester, M., Kriegel, H. P., & Xu, X. (2017). DBSCAN revisited, revisited: Why and how you should (still) use DBSCAN. ACM Transactions on Database Systems, 42(3).

[10] Afshan, A., Kreiman, J., & Alwan, A. (2022). Speaker discrimination performance for "easy" versus "hard" voices in style-matched and -mismatched speech. The Journal of the Acoustical Society of America, 151(2), 1393–1403.

IAFPA conference 2022, 10-13 July 2022, Prague, Czech Republic