

BACKGROUND & OBJECTIVE

- **Proficiency test** recently carried by the *Bundeskriminalamt* (BKA) to evaluate the performance of experts in speaker identification tasks:
 - auditory evaluation
 - Speakers for comparison = **female German twins**
 - widely assumed that twins' voices are similar, and thus recognition of voices is especially difficult (e.g. [1,2]).
- Results:
 - **lack of native knowledge of the language spoken by the twins not a disadvantage for telling the twins apart**
 - informal feedback from participants: **voice quality (VQ) – approached holistically rather than analytically – was the main cue used by non-native listeners to distinguish the twins**
- **Limitations:**
 - limited and idiosyncratic data set (the twins were of advanced age and had lived in different dialectal regions)

SO WHAT NOW?

- BKA test called for the design of a perceptual experiment of a different nature to shed light on how listeners of different L1s perform when assessing the voice of very similar-sounding speakers

In this study we have tested, with a larger twin sample and under controlled conditions of age and dialect, whether the different L1 of listeners affect the perceptual distances between speakers

WHY IDENTICAL TWINS

- **Monozygotic (MZ) twins** are genetically identical and usually share environmental (educational + social) influences
 - Both *organic* (vocal tract anatomy) and *learned* (phonetic choices) variation are minimised in MZ twin pairs.

Importance of twins for voice quality research

Different speakers present isomorphic but not identical vocal tracts, this being one of the shortcomings of perceptual protocols for the assessment of voice quality, such as the VPA:

“Laver’s framework would not be designed to capture the less linguistic aspects of speech, i.e. the relevant ‘phonetic detail’. In other words (...) the small differences in size or shape that two speakers have will make them sound different even if they choose the same articulatory options” [5]

→ Investigations with twins (identical vocal tract) may then prove useful to assess VQ closeness in very similar-sounding speakers.

MATERIALS & METHOD

Subjects: 5 pairs of male MZ twins [6]; native speakers of Standard Peninsular Spanish

- similar age (mean: 21, sd: 3.7)
- similar mean f0 (mean: 113 Hz, sd: 13 Hz)
- similar Euclidean distance (ED) based on perceptual assessment of VQ [7]. EDs produced by ESS [6] via a simplified version of the Vocal Profile Analysis (VPA) scheme [3].

MAJOR SETTING GROUPS											
Key	Labial	Mandib.	Ling. tip	Ling. body	Pharynx	Velo-pharynx	Larynx Height	VT tension	L tension	Phon. Types	
1a	Lip rounding	Close	Advanced	Front & Raised	Constricted	Audible nasal escape	Raised Larynx	Tense	Tense	Falsetto	
1b	Lip spreading	Open	Retracted	Back & Lowered	Expanded	Nasal	Lowered larynx	Lax	Lax	Creak.	
1c	Labiodent.	Protr.								Whisp.	
1d										Harsh.	
1e										Harsh.	

Table 1: Simplified Vocal Profile Analysis Scheme (SVPAS): 10 major setting groups and 26 total settings, with the category key for marking non-neutrality (1a-1e). VT: Vocal Tract; L: Laryngeal; Whisp. = Whispersiness/Breathiness.

	Lab	Mand.	Ling. tip	Ling. body	Pharynx	Velo-pharynx	Larynx Height	VT tension	L tension	Phon. Types	
AGF	0	1a	1a	0	0	0	0	1a	1a	0	
SGF	0	1b	1a	0	0	0	1b	1a	0	1c	0.6
Matchez	1	0	1	1	1	1	0	1	1	1	SMC
AMG	0	1b	1b	0	0	0	0	1a	0	0	
EMG	1a	1b	1a	1b	0	0	1a	1a	0	0	
Matchez	0	1	0	1	1	0	0	1	1	1	SMC
ARJ	0	1a	1a	0	0	0	0	1b	1b	1c	
JRJ	0	1a	0	0	0	0	1b	1b	1b	1c	
Matchez	1	0	1	1	1	1	0	1	1	1	SMC
ASM	1a	0	0	0	1b	0	1b	0	0	1c	
RSM	1a	1a	0	0	1b	0	1b	1a	1b	0	
Matchez	1	0	1	1	1	1	0	0	0	0	SMC
DCT	1a	0	0	1a	0	0	1a	1a	1a	1d	
JCT	0	0	1a	0	1a	1b	1a	1a	1a	1d	
Matchez	1	1	0	0	0	1	1	1	1	1	SMC

Table 2: Summary of SMCs for all twin pairs. Mean SMC: 0.68, indicating that around 7 VQ settings were shared on average by the twin pairs.

Stimuli: Speech samples (~3 secs) extracted from semi-directed spontaneous conversations [6]. Declarative sentences of diverse neutral topics.

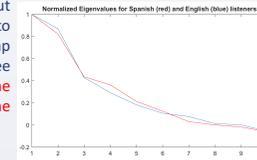
Listeners: Group 1: native Spanish speakers (N=20; age range 22-51, mean 33). Group 2: native English speakers with no knowledge of Spanish (N=20; age range 19-35, mean 25).

Design of perceptual test: Multiple Forced Choice experiment. 90 different-speaker pairings, i.e. each speaker compared with everyone else. Stimuli presented in random order. Listeners indicated degree of similarity of pair on 1–5 scale. They didn't know that the stimuli included twin pairs. The test was run on a PC with HQ headphones. Short pre-test for familiarization.

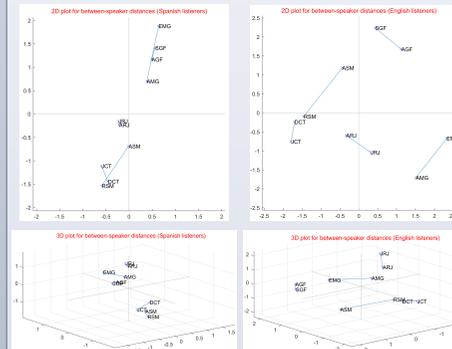
Analysis method: Degree of perceived similarity measured using Multidimensional Scaling (MDS), a means of visualizing the level of similarity of individual cases in a dataset and of detecting meaningful underlying dimensions that explain observed similarities or dissimilarities (distances) [4].

RESULTS

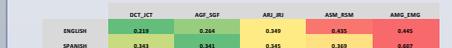
MDS analyses carried out using similarity scores, to construct a perceptual map of the speakers. The scree plot (right) shows the relative magnitude of the sorted Eigenvalues.



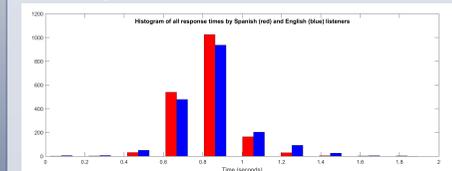
7 dimensions necessary to accurately reproduce between-speaker distances in the perceptual space, but MDS results typically visualized using only the first 2 or 3 dimensions.



Normalized intra-pair Euclidean distances (7 dimensions):



Reaction times very similar for Spanish (mean: 0.82 secs; std: 0.14) and English listeners (mean: 0.84; std: 0.18).



DISCUSSION

↑ All speakers are closer in the perceptual space. Does this imply that knowledge of the linguistic content make the task more difficult? Distraction effect of the message?

↑ Better detection of very similar speakers (i.e. twins). Smaller distances between these in comparison with English listeners. Note also the different magnitude of the plots.

↑ Most similar twin pair: AGF-SGF and ARJ-JRJ. Most different twin pair: AMG-EMG.

↑ All speakers are more spread in the perceptual space. Some twin pairs are very far apart, which even makes them have an unrelated speaker as their closest speaker.

↑ Most similar twin pair: DCT-JCT. Most different: AMG-EMG.

CONCLUSIONS

- Eigen-decomposition: 7 main dimensions explain similarity decisions by listeners (both English and Spanish)
 - voice is highly multidimensional and that reducing perceptual space to 2D or 3D may be misleading
 - similarity of the relative magnitude of the sorted Eigenvalues suggest that similar perceptual strategies operate for both listener groups.
- Almost the same ranking of twin similarity for both listener groups could indicate the same cue prominence, i.e. regardless of familiarity with the language spoken or understanding of the linguistic content, both groups show ~ same ranking of similarity of twin pairs. **Exception:** AGF-SGF most similar for Spaniards (VQ analysis: tense VT & advanced tongue tip) while DCT-JCT most similar for English (VQ analysis: harshness & raised larynx). Different perceptual salience of VQ settings?
- Equivalent reaction times point to similar listening strategies ('gut' impressions; holistic VQ perception).
 - However, qualitative feedback from participants also point to other cues: mainly rhythmic aspects but also segmental features.
- Besides: your twin may not necessarily be your best impostor (e.g. RSM: closer perceptual distance with unrelated speaker)
- **Future work:** Correlate Euclidean distances obtained from VQ holistic perception with component-featural analysis of VQ
- **Other future work:** (1) Sort listeners in musical and non-musical training; (2) Test with listeners of other languages (e.g. Germans - would these obtain similar results as English?)

REFERENCES

- [1] Decoster, W., Van Gysel, A., Vercammen, J. & Debruyne, F. (2000) Voice similarity in identical twins. *Acta Oto-Rhino-Laryngologica Belgica*, 55: 49-55.
- [2] Künzle, H.J. (2010) Automatic speaker recognition of identical twins. *Int. J. Speech, Language & the Law* 17: 251-277.
- [3] Laver, J. (1980) *The phonetic description of voice quality*. Cambridge: CUP.
- [4] McDougall, K. (2013) Assessing perceived voice similarity using Multidimensional Scaling for the construction of voice parades. *Int. J. Speech, Language & the Law* 20: 163-172.
- [5] Nolan F. (2005) Forensic speaker identification and the phonetic description of voice quality. In W.J. Hardcastle & J. Beck (eds.) *A Figure of Speech: A Festschrift for John Laver*. 385–411. Mahwah, NJ: Erlbaum.
- [6] San Segundo, E. (2014) *Forensic speaker comparison of Spanish twins and non-twin siblings*. Doctoral dissertation, Menéndez Pelayo Int. Univ.
- [7] San Segundo, E. & Mompesán, J.A. (2016) Voice quality similarity based on a simplified version of the Vocal Profile Analysis. *Sociolinguistics Symposium 21*, University of Murcia, Spain, 15-18 June.

ACKNOWLEDGEMENTS

Thanks to Andrew MacFarlane and Duncan Robertson for help with measuring reaction time. A special thanks to Dr. Juana Gil for letting ESS run the test with Spanish listeners at the Phonetics Lab (CSIC, Madrid).