



UNIVERSITÀ  
DEGLI STUDI  
FIRENZE  
**DINFO**  
DIPARTIMENTO DI  
INGEGNERIA  
DELL'INFORMAZIONE

8th  
INTERNATIONAL  
WORKSHOP

MODELS AND  
ANALYSIS  
OF VOCAL  
EMISSIONS  
FOR  
BIOMEDICAL  
APPLICATIONS

December 16-18, 2013  
Firenze, Italy



**PROCEEDINGS**



**Reference:** San Segundo, E., & Gómez-Vilda, P. (2013). Voice biometrical match of twin and non-twin siblings. In *Proceedings of the 8th International Workshop Models and analysis of vocal emissions for biomedical applications, Firenze, Italy* (pp. 253--256).

# VOICE BIOMETRICAL MATCH OF TWIN AND NON-TWIN SIBLINGS

Eugenia SanSegundo<sup>1</sup>, Pedro Gómez-Vilda<sup>2</sup>

<sup>1</sup>Phonetics Lab., Institute of Language, Literature and Anthropology, Spanish National Research Council (CSIC)  
C/ Albasanz 26-28, 28037 Madrid, Spain

<sup>2</sup>NeuVox Laboratory, Center for Biomedical Technology, Universidad Politécnica de Madrid, Campus de Montegancedo, s/n, 28223 Pozuelo de Alarcón, Madrid, Spain  
e-mails: [eugenia.sansegundo@cchs.csic.es](mailto:eugenia.sansegundo@cchs.csic.es), [pedro@fi.upm.es](mailto:pedro@fi.upm.es)

**Abstract:** The similarity in twins' voices has always been an intriguing issue in forensic speaker matching, and has become an important research matter recently. The present work is a preliminary study of exploratory character diving into the similarities of monozygotic (MZ) and dizygotic (DZ) twins' phonation under the point of view of vocal fold biomechanics. The study extends to other siblings' and unrelated speakers' phonation. Estimates of biomechanical parameters obtained from vowel fillers are used to produce bilateral matches between MZ and DZ twins and siblings, and unrelated speakers. These results show interesting relationships regarding genetic load and ambient factors in the adoption of phonation styles.

**Keywords:** voice production, forensic pattern matching, phonation styles, glottal source features, twins' voice.

## I. INTRODUCTION

Recent studies in voice quality are conducted towards the evaluation of phonation performance in relation to either professional voice care, or in meta-acoustic knowledge (neurological deterioration, emotion detection, forensic applications, etc.) These fields of study are becoming more and more demanded nowadays. The aim of the present work is to study the similarities and differences of phonation characteristics in twins' voices, including monozygotic (MZ) as well as dizygotic (DZ) twins for specific forensic use, not disregarding other fields of application, as the clinical one, although this is not the main aim of the paper. A reference to previous work on twin voice quality analysis and vocal performance of interest is that of Van Lierde et al. [1]. The quality measurements used were perceptual GRBAS, breathing performance, fundamental frequency, jitter and shimmer, and the Dysphonia Severity Index. However, the study focused only on monozygotic siblings (MZ). Another relevant reference is that of Cielo et al. [2], although the twin sample used was quite small (2 MZ pairs, one per gender). Their analysis is interesting as far as they use some features not been considered in twins' voice studies before, namely vocal onset and harmonic characterization. The work of Fuchs et al. [3] found that

the voices of MZ twins showed more similarity among themselves than those of non-similar speakers regarding vocal range, highest and lowest fundamental frequency, prosodic pitch line, maximum intensity, number of overtones and intensity vibrato.

The study of twins' voices can be approached from many perspectives. Stemming from a typical phonetic division, they may be classified into perception, acoustics or articulation. Some of the acoustic-related studies dealing with voice-quality or glottal parameters have been reviewed in [[4]]. Since perceptual or articulation-based approaches are less relevant for the purpose of this work, we will consider those studying twins' voices from an automatic perspective. The system by Scheffer et al. [[5]] was able to identify twins with a good performance (85% of correct identifications) using MFCC (Mel Frequency Cepstrum Coefficients). The residual error (speakers who were not correctly detected as twins of their actual twins) would suggest that "the twin of a speaker is not necessarily the most difficult impostor for an automatic speaker recognition system" ([[5]]: 2). The automatic system by Ariyaeinia et al. [[6]] used LPCC (Linear Predictive Coding-Derived Cepstral) parameters, and the speaker representation was based on adapted Gaussian Mixture Models (GMMs). The results showed that the use of long test utterances led to smaller error rates than short ones. Both KyungWha [[7]] and Künzel [[8]] used *Batvox*; the former studied Korean female twin pairs (17 MZ, including 1 triplet and 5 DZ) and the latter studied German male and female twin pairs. The results in [[7]] showed that every twin speaker was correctly identified in the same speaking style condition (reading speech). The performance of the system in [[8]] was better for male than for female voices.

The present work focuses on studying phonation marks (including biomechanical parameters) of relevance in the biometrical description of phonation [[10], [11]]. The working hypothesis is that phonation cycle quotients and biomechanics may offer differentiation capabilities among MZ, DZ and control speakers not explored already. The paper is organized as follows: A description of the materials and methods used in the study is given in section 2. In section 3 results obtained from the bilateral tests and matches of 16 male speakers are discussed. Conclusions are presented in section 4.

## II. METHODS

Recordings from 40 male native speakers of Spanish (spontaneous conversation) were taken at a sampling rate of 44,100 Hz and 16 bits using HQ microphones in an isolated room. The distribution of speakers was: 7 MZ pairs, 5 DZ pairs, 4 pairs of non-twin siblings and 4 pairs of controls (non-relatives). Spontaneous fillers (long [ε] vowels maintained during more than 200 ms produced by speakers in words like “que”, “de”, or in hesitation marks like “eeh...” etc.) were used in the study. Recordings from two sessions separated by a 3-week interval were taken per speaker. Speech recordings were around 10 min long, an average of 8-10 fillers found in each recording.

A set of biomechanical parameters as body and cover dynamic mass and stiffness was estimated from the glottal source by inverse filtering [9]. The inter-cycle unbalances of these parameters were also used. Open, Close and Return Quotients were added to the parameter set as well as Contact Gap Defects. The parameter set was completed with jitter, shimmer and NHR ratio to produce a feature vector of 65 parameters given as  $x_{sij}$ , where  $s$  refers to the subject,  $i$  is for the session, and  $j$  for the filler. Pair-wise parameter matching experiments were carried out by likelihood ratio contrasts used in forensic voice matching [11]. The test is based on two-hypotheses contrasts: that the conditional probability between voice samples  $Z_a=\{x_{aij}\}$  and  $Z_b=\{x_{bij}\}$  (from the two subjects under test,  $a$  and  $b$ ) is larger than the conditional probability of each subject relative to a Reference Speaker's Model  $\Gamma_R$  in terms of logarithmic likelihood

$$\lambda_{ab} = \log \left[ \frac{p(Z_b | \Gamma_a)}{\sqrt{p(Z_a | \Gamma_R)p(Z_b | \Gamma_R)}} \right] \quad (1)$$

where conditional probabilities have been evaluated using Gaussian Mixture Models ( $\Gamma_a, \Gamma_b, \Gamma_R$ ) as

$$\begin{aligned} p(Z_b | \Gamma_a) &= \Gamma_a(Z_b); \\ p(Z_a | \Gamma_R) &= \Gamma_R(Z_a); \\ p(Z_b | \Gamma_R) &= \Gamma_R(Z_b) \end{aligned} \quad (2)$$

Following this background, the Forensic Voice Evidence Evaluation Framework is a two-step process:

- **Step 1. Model Generation.** A model representative of the normative population set considered (male subjects between 18-52 years-old) was created on recordings  $Z_R=\{x_{Rjk}\}$ , as a Gaussian Mixture Model  $\Gamma_R=\{\mathbf{w}_R, \boldsymbol{\mu}_R, \mathbf{C}_R\}$ ,  $\mathbf{w}_R, \boldsymbol{\mu}_R$  and  $\mathbf{C}_R$  being the set of weights, averages and covariance matrices associated to each Gaussian Probability Distribution in the set.
- **Step 2. Score Evaluation.** The material under evaluation will be composed of different parameterized voice samples in matrix form  $Z_a=\{x_{aj}\}$ , where  $1 \leq j \leq J_a$  is the sample index, each sample being a vector  $\mathbf{x}_{aj}=\{x_{aj1} \dots x_{ajM}\}$  from vowel-like segments conveniently parameterized. Similarly, the set of the correspondent speaker to be matched

will be given as  $Z_b=\{x_{bj}\}$ , where  $1 \leq j \leq J_b$  will be the sample index, each sample being a vector  $\mathbf{x}_{bj}=\{x_{bj1} \dots x_{bjM}\}$ .

The conditioned probability of a sample from speaker  $a$  matching speaker  $b$  will be estimated as

$$\Pr(\mathbf{x}_{bj} | \Gamma_a) = \frac{1}{(2\pi)^{M/2} |\mathbf{C}_a|^Q} e^{-1/2(\mathbf{x}_{bj}-\boldsymbol{\mu}_a)^T \mathbf{C}_a^{-1}(\mathbf{x}_{bj}-\boldsymbol{\mu}_a)} \quad (3)$$

Similarly the conditioned probability of a sample from speaker  $a$  matching the Reference Model will be

$$\Pr(\mathbf{x}_{aj} | \Gamma_R) = \frac{1}{(2\pi)^{M/2} |\mathbf{C}_R|^Q} e^{-1/2(\mathbf{x}_{aj}-\boldsymbol{\mu}_R)^T \mathbf{C}_R^{-1}(\mathbf{x}_{aj}-\boldsymbol{\mu}_R)} \quad (4)$$

Finally the conditioned probability of a sample from speaker  $b$  matching the Reference Model will be

$$\Pr(\mathbf{x}_{bj} | \Gamma_R) = \frac{1}{(2\pi)^{M/2} |\mathbf{C}_R|^Q} e^{-1/2(\mathbf{x}_{bj}-\boldsymbol{\mu}_R)^T \mathbf{C}_R^{-1}(\mathbf{x}_{bj}-\boldsymbol{\mu}_R)} \quad (5)$$

A full description of this methodology is given in [[12]].

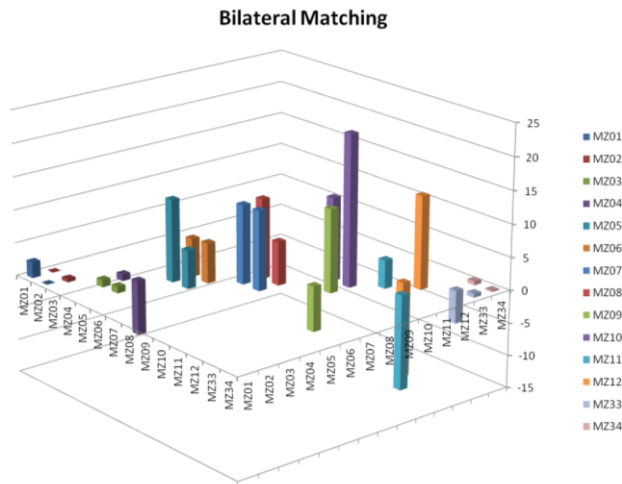
## III. RESULTS AND DISCUSSION

The composition of the sample was the following: 14 subjects are MZ siblings in 7 pairs (numbered as 01-02, 03-04, 05-06, 07-08, 09-10, 11-12 and 33-34), 10 subjects are DZ siblings in 5 pairs (corresponding to speakers numbered as 13-14, 15-16, 17-18, 19-29 and 45-46), 8 subjects are non-twin brothers (RS) in 4 pairs (numbered as 21-22, 23-24, 47-48 and 49-50) and 8 subjects are not known to have any familiar relationship (US), grouped also as 4 pairs (25-26, 27-28, 29-30 and 31-32). Speakers were matched in: a) different-session intra-speaker tests (I: intra-speakers), b) inter-speaker tests (O: inter-speakers). A priori expectations assume that MZ should show the largest LLRs, followed by DZ, then by non-twin siblings; non-related speakers are expected to show the lowest LLRs. The baseline is defined by a reference background set composed of 20 speakers (set B). Scores are qualified as Strong Likeness if above 1, Weak Likeness if between 1 and -1 and Unlikeness if below -1. The hypotheses tested were the following:

- H1. Intra-speaker tests should show large LLRs.
- H2. MZ inter-speaker tests should show large LLRs.
- H3. DZ inter-speaker tests should show also large LLRs although not that large as H1 or H2.
- H4. RS inter-speaker tests should show LLRs at least over the background baseline (fixed at  $\lambda = -10$ ).
- H5. US inter-speaker tests should show LLR's aligned with the background baseline.

The results of the matching tests are summarized in Table 1 (see end of paper). The results contradicting the strongest hypotheses (H1 and H2) are marked in bold. Four speakers out of the total of 40 appear to be in the limit of H1 (03, 48, 49 and 50), five others show strong intra-speaker dissimilarity (04, 09, 15, 20 and 33), and one shows very strong self-dissimilarity (25), therefore

10 out of 40 do not fulfil H1. The rest of the speakers show weak or strong self-similarity in inter-session tests, fulfilling H1. Regarding H2 we find only one out of seven pairs not fulfilling it (11 vs 12). Hypothesis 3 is not fulfilled in one out of five pairs (17 vs 18). H4 is fulfilled in all four cases. Only one pair of unrelated subjects is slightly over the baseline (27 vs 28) out of 4 cases fulfilling H5. The cases affecting only to MZ siblings have been depicted in Fig. 1 for special discussion.



**Fig. 1** Summary of the results for the MZ tests.

The 3 intra-speaker tests out of 14 which do not fulfil H1 correspond to relatively large negative column values (04, 09 and 33) as well as one inter-speaker test not fulfilling H2 (11 vs 12). Two twin pairs show good fulfilment of H1 and H2 (05, 06, 07 and 08), another twin pair do show a weak fulfilment of H1 and H2 (01 and 02), two twin pairs show weak fulfilment of H2, and irregular fulfilment of H1 (03, 04, 33 and 34). Another twin pair shows strong fulfilment of H2 and irregular fulfilment of H1 (09 and 10) and another pair shows good fulfilment of H1 but irregular fulfilment of H2 (11 and 12). Some words have to be said about intra-speaker fulfilment of H1: it is unclear why 10 out of 40 speakers do show self-unlikeness in a larger or smaller extent when one session phonation is tested against another. Several reasons have been considered, as changes in phonation due to emotional stress or even temporary pathological conditions. Excluding weak self-unlikeness the number of cases would be 6 out of 40, which is still a large figure. Possibly some normalization on the selection of the speaker's most characteristic phonation patterns could help in reducing this apparently large value. Regarding H2 the number of non-fulfilments seems smaller (1 out of 7 pairs). Reasons for dissimilarities in MZ within-pair comparisons seem somehow different. The most plausible reason that we can pinpoint is the nature-nurture dichotomy: in other words, the behavioural component of phonation as opposed to genetic reasons

(phonation characteristics may be due to learned styles as much as to biological imprinted patterns).

#### IV. CONCLUSIONS

The results of the study show some interesting considerations. Regarding H1 it seems that there are certain speakers who do not show strong intra-speaker similarity (6 out of 40 are in this situation). The immediate reflection is if these could be labelled as "goats" in Doddington's Zoo [[13]]. As far as H2 is concerned it seems that most MZ twins show reasonable inter-speaker (within-pair) similarity except in one pair out of 7. Whether this could be due to behavioural rather than to genetic factors is an open question. In DZ twins (H3) the situation is similar (only 1 out of 5 pairs show low inter-speaker scores). Non-twin brothers fulfil H4 relatively well, since all 4 pairs considered showed scores over the background baseline. Finally non-relative subjects showed scores well around the background baseline giving a good description of what would be considered the normal situation in unrelated speakers. A possible complementary explanation involves the 65 parameter set in such comparisons where some of them may show a greater influence from both genetic and environmental factors. If only the comparisons of MZ twin pairs had yielded large matches, the only explanation possible would be genetic influence. However, the fact that similar values are obtained for MZ and DZ twins cannot lead to that conclusion. The impact of external factors (like a similar living and educational environment, same age, etc.) may be more relevant than it may be thought a priori in this kind of voice studies. Further research would be necessary in order to study the role of each specific parameter intervening in the results, and to extend the study to more speakers.

**Acknowledgments:** This work is supported by an FPU grant from the Ministry of Education, a grant from the International Association for Forensic Phonetics and Acoustics, and by grants TEC2009-14123-C04-03 and TEC2012-38630-C04-04 from *Plan Nacional de I+D+i*, Ministry of Economy and Competitiveness of Spain.

#### REFERENCES

- [1] Van Lierde, K. M., Vinck, B., De Ley, S., Clement, G., and Van Cauwenberge, P. "Genetics of vocal quality characteristics in monozygotic twins: a multiparameter approach", *Journal of Voice*, Vol. 19, No. 4, 2005, pp. 511-518.
- [2] Cielo, C. A., Agustini, R. and Finger, L. S., "Características vocais de gêmeos monozigóticos", *Revista CEFAC*, Vol. 14, No 6, 2012, pp. 1234-1241 (in Portuguese, summary in English).
- [3] Fuchs, M., Oeken, J., Hotopp, T., Täschner, R., Hentschel, B. and Behrendt, W., "Die Ähnlichkeit monozygoter Zwillinge hinsichtlich

- [4] San Segundo, E. and Gómez-Vilda, P., “Voice Biometrical Match of Twin and non-Twin Siblings”, *Proc of the 1st Multidisciplinary Conf. of Users of Voice, Speech and Singing*, Las Palmas de Gran Canaria, 27-28 June 2013, pp. 132-136.
- [5] Scheffer, N., Bonastre, J-F., Ghio, A. and Teston, B. “Gémellité et reconnaissance automatique du locuteur”, *Actes, Journées d’Etude sur la Parole (JEP)*, 2004, pp. 445-448.
- [6] Ariyaeinia, A., Morrison, C., Malegaonkar, A. and Black, S., “A test of the effectiveness of speaker verification for differentiating between identical twins”, *Science and Justice*, Vol. 48, 2008, pp. 182-186.
- [7] KyungWha, K., “Automatic speaker identification of Korean female twins”, *Proc. 19th Ann. Conf. of the Int. Assoc. for Forensic Phon. and Acoust. (IAFPA)*, Trier, Germany, 18-21 July 2010, p. 21.
- [8] Künzel, H., “Automatic speaker recognition of identical twins”, *Int. Journal of Speech Language and the Law*, Vol. 17, No. 2, 2010, pp. 251-277.
- [9] Gómez, P., Fernández, R., Rodellar, V., Nieto, V., Álvarez, A., Mazaira, L. M., Martínez, R. and Godino, J. I., “Glottal Source Biometrical Signature

- [10] Gómez, P., Rodellar, V., Nieto, V., Martínez, R., Álvarez, A., Scola, B., Ramírez, C., Poletti, D., and Fernández, M., “BioMet@Phon: A System to Monitor Phonation Quality in the Clinics”, *Proc. eTELEMED 2013: The Fifth Int. Conf. on e-Health, Telemedicine and Social Medicine*, Nice, France, 2013, pp. 253-258.
- [11] González, J., Rose, P., Ramos, D., Toledano, D. T. and Ortega, J., “Emulating DNA: Rigurous Quantification of Evidential Weight in Transparent and Testable Forensic Speaker Recognition”, *IEEE Trans. On Audio, Speech and Lang. Proc.*, Vol. 15, No. 7, 2007, pp. 2104-2115.
- [12] Gómez, P., Mazaira, L. M., Hierro, J. A. and Nieto, R., “Distance Metrics in Voice Forensic Evidence Evaluation using Dysphonia-relevant Parameters”, *Proc of the VI Jornadas de Reconocimiento Biométrico de Personas*, Las Palmas de Gran Canaria, January 26-27, 2012, pp. 169-178.
- [13] Doddington, G., Liggett, W., Martin, A., Przybocki, M., & Reynolds, D., “Sheep, goats, lambs and wolves: A statistical analysis of speaker performance in the NIST 1998 speaker recognition evaluation”, NIST, Gaithersburg, MD, 1998.

[illegible]