

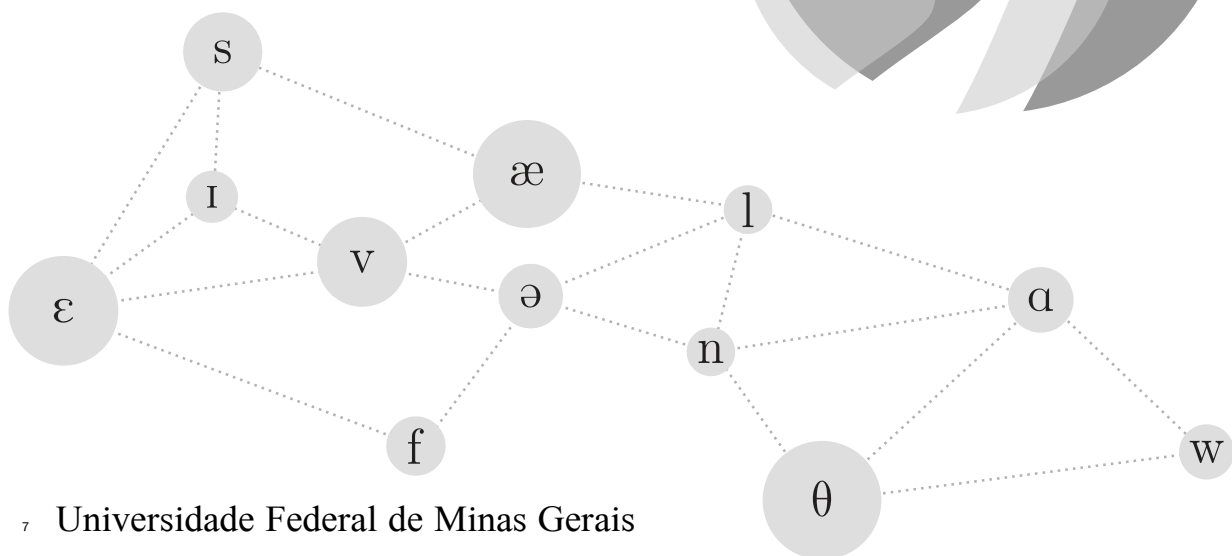
<sup>1</sup> Thaís Cristófato-Silva, Hani Yehia, Leonardo Araujo,  
<sup>2</sup> Maria Cantoni, Magnum Madruga and Adriano Vilela

# <sup>3</sup> EICEFALA 2021

<sup>4</sup> International Meeting on Speech Sciences

<sup>5</sup> Advances in speech and L2 processing

<sup>6</sup> SEVENTH EDITION



<sup>7</sup> Universidade Federal de Minas Gerais

9 Copyright © 2021 Thaís Cristófato-Silva, Hani Yehia, Leonardo Araujo,

10 Maria Cantoni, Magnum Madruga and Adriano Vilela

11 PUBLISHED BY UNIVERSIDADE FEDERAL DE MINAS GERAIS

12 Licensed under the Apache License, Version 2.0 (the “License”); you may not use this file  
13 except in compliance with the License. You may obtain a copy of the License at [http:](http://www.apache.org/licenses/LICENSE-2.0)  
14 [//www.apache.org/licenses/LICENSE-2.0](http://www.apache.org/licenses/LICENSE-2.0). Unless required by applicable law or agreed  
15 to in writing, software distributed under the License is distributed on an “AS IS” BASIS,  
16 WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.  
17 See the License for the specific language governing permissions and limitations under the  
18 License.

19 First printing, October 2021

## 20

## 21

22

25

26

28

29

30

31

33

## 34

35

38

39

41

Production of English [Cs] clusters by Brazilian speakers:  
effects of orthography, phonological environment and task  
type 46

WELLINGTON ARAUJO MENDES JUNIOR

Radial Basis Function Artificial Neural Network for Au-  
tomatic Identification of Interlanguage Transfer Phenom-  
ena 58

ATOS APOLLO SILVA BORGES, BRUNO FERREIRA DE SOUSA, ARATUZA RO-  
DRIGUES SILVA ROCHA, WILSON JÚNIOR DE ARAÚJO CARVALHO, FÁBIO ROCHA  
BARBOSA, RONALDO MANGUEIRA LIMA JÚNIOR

PART III METHODS FOR SPEECH DATA COLLEC-  
TION, PROCESSING AND ANALYSIS

Behavioral and Neurophysiological Representations of Speech  
Phonemic Units 68

ADRIELLE C. SANTANA, ADRIANO V. BARBOSA, HANI C. YEHIA, RAFAEL  
LABOISSIÈRE

PART IV BRAZILIAN PORTUGUESE PHONETICS  
AND PHONOLOGY

The implementation of phonic voicing contrast in children's  
speech: some explorations of clinical data 79

FABIANA NOGUEIRA GREGIO, ZULEICA CAMARGO

PART V TEST

# 61 Introduction

62 EICEFALA is an event promoted by UFMG laboratories from Facul-  
63 dade de Letras (Laboratório de Fonologia) and Escola de Engenharia  
64 (CEFALA).

65 The comprehension of human communication involving speech  
66 acoustics and gestures requires knowledge from various fields of sci-  
67 ence, such as phonetics, phonology, linguistics, acoustics, mechanics,  
68 mathematics, physiology, neuroscience and computer science. The  
69 objective of the 7th EICEFALA is to present and discuss theoretical  
70 and methodological techniques to researchers from the several ar-  
71 eas of knowledge working with speech science: linguists, engineers,  
72 physicists, speech therapists, musicians, etc. It is expected that the par-  
73 ticipants of the event find, at EICEFALA, a transdisciplinary forum to  
74 address questions related to spoken communication.

<sup>75</sup> PART I:

<sup>76</sup> PLENARY SPEAKERS

Fairy tales are more than true: not because  
they tell us that dragons exist, but because  
they tell us dragons can be beaten.

C.K. CHESTERTON

# Using statistical learning techniques to determine Cantonese lexical tones from the acoustic and visual components of speech

JOÃO VITOR POSSAMAI DE MENEZES<sup>1</sup>, HANI CAMILLE YEHIA<sup>1</sup>,  
ADRIANO VILELA BARBOSA<sup>1</sup>

<sup>1</sup> . Universidade Federal de Minas Gerais

This mini-course presents an introduction to the use of statistical learning techniques to speech processing problems. More specifically, we show how classification techniques can be used to predict lexical tones in Cantonese from the associated measurements of both the acoustic and the visual (to a lesser degree) components of speech. The acoustic and visual data we use were recorded during a speech production experiment where a native speaker of Cantonese produced a set of words spanning the full range of Cantonese tones. The visual data consists of 3D trajectories of markers on the subject's face and head recorded with an Optotrak. The acoustic component is represented by F0 trajectories extracted from the speech acoustics. The idea is to use the F0 and marker trajectories as input vectors to train classifiers to predict the lexical tones. However, these trajectories cannot be used directly because they have different durations for different tokens (utterances), whereas all input vectors to the classifiers must have the same dimension. In order to make all input vectors the same length, regardless of the duration of the utterances, all trajectories (both F0 and markers) are approximated by polynomials of a given order and represented by the corresponding coefficients. The polynomial coefficients are then used as input vectors to train different classification models (LDA, SVM, K-nearest neighbors, etc). The performance of the models is estimated

104 by means of k-fold cross validation. Although the statistical learning  
105 techniques we present are applied to a specific problem (estimating  
106 Cantonese lexical tones from the acoustical and visual components of  
107 speech), they are general and can be equally applied to a wide range of  
108 problems. All procedures presented in the mini-course are developed in  
109 the R language.



# The multiple dimensions of speech: old questions and new challenges

DIDIER DEMOLIN<sup>1</sup>

1 . Laboratoire de Phonétique et Phonologie

CNRS-UMR 7018

Human speech, a product of the evolution of primates, can in essence be defined in terms of a signal. This is an acoustic wave varying over time with amplitude and frequency modulations, due to the articulatory movements of the vocal tract's organs. To perform these movements, motor controls are required, whose interactions with the aerodynamic parameters produce the acoustic signal. The main objective of research in this domain is to understand which primary principles, biological, physical and cognitive, to be based on to explain the production and perception of speech in the world's languages and to make the fundamental question: how does it work?

Among the main fields of activity involved in the study of sounds and sound systems of languages are the engineering sciences with the dimensions of automatic processing (speech recognition and synthesis); phonetics and phonology (the linguistic aspects); and pathological aspects (how to explain what doesn't work anymore or less well). This includes knowledge of similar fundamental principles. To these dimensions a readded physics, biology, cognition and neuroscience. These fields involves in-depth knowledge of various interconnected fields to explain how sounds and sound systems work. Therefore in addition to the symbolic dimension, anatomical, physiological, acoustic, aerodynamic, articulatory, auditory, proprioceptive, historical (phylogenies and diachrony), ecological, temporal, dynamic and self-organized as-

pects can, and should, be integrated in the explanation of the studied phenomena.

The complexity and interactions of these dimensions find new light in the paradigms resulting from the study of complex systems, which makes it possible to address old issues again, such as the search for a possible speech code, invariants and primitives. From these issues, others arise, such as the understanding of the open or closed nature of sound systems, which is far from being resolved. Explaining the diversity, complexity and dynamics of sound systems involves understanding the nature of variation in speech phenomena. How can we show that spontaneous speech, laboratory speech and pathological aspects are based on the same principles?

The evolution of theory, models, new statistical tools, computational, big data and deep learning tools, allow these issues to be addressed in a new light. New measuring instruments such as real-time magnetic resonance imaging, functional magnetic resonance, three-dimensional or four-dimensional ultrasound, digital endoscopy, electroencephalography (EEG) and many other recent tools make it possible to accurately observe, measure and quantify speech phenomena as well as bring to discussion fundamental issues still unresolved or poorly understood.

The lecture will discuss the controlled and automatic aspects involved in the control of breathing in speech, issues in speech embodiment, the quantal aspects of speech, the importance of thresholds values in aerodynamic and acoustic parameters, types of feedback (acoustic and proprioceptive) in speech phenomena and new ways to explain and formalize the source, the initiation and propagation of sound changes. This last point by using and adapting population ecology models to speech.

# Some questions on L2 speech as related to colonialism

ELEONORA ALBANO<sup>1</sup>, ANTONIO PESSOTTI<sup>1</sup>, CARLA DIAZ<sup>1</sup>

1. LAFAPE-IEL-UNICAMP & CNPq

The first aim of this talk is to revisit the question of phonetic drift in  $L_2$  speech in light of new data and theory. The new data consist of a sizeable set of acoustic-phonetic measures of the speech of Quechua and Spanish monolinguals and bilinguals residing in Peru. The theoretical innovation draws on two sources: the relatively familiar concept of accommodation, introduced by Giles et al. in sociolinguistics, and the less familiar concept of coloniality, introduced by Quijano (2000) in sociology. At the same time, it aims at showing that phonetic analysis based on gestural phonology can open new avenues for exploring the relationship between these two concepts as explanans for  $L_2$  pronunciation in an ethnically diverse environment.

Accommodation refers to a “constant movement toward and away from others, by changing one’s communicative behavior” (GILES & OGAY, 2007). It encompasses speech and various other communicative behaviors. Moreover, it has a convergent side – enhancing similarities between interlocutors – and a divergent one – enhancing differences between interlocutors. Both can occur between two or more people or within and across speech communities. Some acoustic phonetic parameters have been useful to tap such shifts (e.g., VOT, as in the pioneering work of SANCIER & FOWLER, 1997).

The concept of coloniality refers to “how colonial patterns of power and inequality exceed the spatial and temporal boundaries of empire and colony” (ROCHE, 2019). It aims at dealing with the epistemology

underlying the pervasive replication of colonial social, economic, and cultural practices in postcolonial societies (QUIJANO, 2000).

We will start by revisiting earlier work on phonetic drift in  $L_2$  conducted at our lab – Laboratório de Fonética e Psicolinguística (LAFAPE). Ramirez et al. (2011) showed that contact situations may exhibit intralinguistic phonetic drift in both  $L_1$  and  $L_2$ . In turn, Albano et al. (2020) reported preliminary observations of intralinguistic drift attributable to language attrition in Quechua/Spanish bilinguals residing in Brazil.

We believe that the understanding of the results of both of these works can be considerably improved by reference to the above-defined concepts. In particular, some intriguing signs of partial loss of Quechua stop distinctions shown by the expatriated Peruvians can be interpreted as mistiming of articulatory gestures converging toward those of the two hegemonic languages (namely, Spanish and Portuguese).

Then we will move on to inquire how the study conducted in Peru can elucidate our questions about Quechua/Spanish relations. All data collection on this topic was part of Carla Diaz's requisites for completing her bachelor and master's degrees in linguistics (DIAZ, 2018; 2021).

Carla recorded 10 Spanish monolinguals and 10 Quechua/Spanish bilinguals in Lima in August 2019. Then she travelled to Cuzco to record 11 monolingual Quechua speakers, with the help of a bilingual friend specializing in Quechuan literature.

The corpus, similar to that of Albano et al. (2020), focused on Quechua and Spanish stop contrasts. The analysis, likewise, employed measurements that have been used in the description of Quechua: VOT, amplitude of the stop burst,  $f_0$ , and  $H1 - H2$ .

The results show that, unlike the residents of Brazil, the residents of Lima have no trouble distinguishing stops within the Quechua series or differentiating them from Spanish stops. The remarkable fact is that divergence from Spanish was more frequent than convergence. Moreover, certain distinctions were enhanced by shifting the acoustic parameters beyond the values of the monolingual group.

After Carla defends her thesis, we are planning to conduct finer-grained analyses considering the linguistic values and attitudes captured by our sociolinguistic questionnaire. May we succeed in helping unravel the coloniality issues behind the subtle attempts of Peruvian Quechua speakers at resisting diglossia and language loss.

231 Bibliography

232 Howard Giles, Donald M. Taylor, and Richard Bourhis. Towards a  
233 theory of interpersonal accommodation through language: some  
234 canadian data. *Language in Society*, 2(2):177–192, October 1973.  
235 DOI: 10.1017/s0047404500000701. URL [https://doi.org/10.](https://doi.org/10.1017/s0047404500000701)  
236 1017/s0047404500000701.

237

238

# How to model the influence of orthography on L2 representations with BiPhon Neural Networks

239

SILKE HAMANN<sup>1</sup>, CHAO ZHOU<sup>1</sup>

240

1 . University of Amsterdam and University of Lisbon

241

Many studies have shown that written forms influence the acquisition of a second language. This influence can be helpful, as is the case of the English /æ/-/ɛ/ contrast that is notoriously difficult for Dutch learners but where the written form can aid in the creation of the distinction [?Escudero and Wanrooij, 2010]. But orthography can also cause the creation of so-called ghost contrasts, which do not exist in the L2, as is the case with the intervocalic singleton/geminate contrast in the L2 English of Italian speakers [Bassetti, 2017, Hamann, 2018].

249

In this talk, we illustrate how such orthographic influences on the creation of L2 representations can be formalized, by this yielding theoretical predictions that can be tested again in experimental studies. Our formalization is performed with a symbolic neural network based on the Bidirectional Phonetics-Phonology model [Boersma, 2007] and its extension by a reading grammar [Hamann and Colombo, 2017].

255

Our main data comes from an experimental study on Mandarin [Zhou and Hamann, 2020]: 23 L1-Mandarin speakers with no prior knowledge of EP (naïve listeners), representing the initial stage of L2 acquisition, performed a delayed-imitation task. They were presented with EP nonce words containing /r/ in intervocalic onset (e.g., parafa) or word-internal coda (e.g., parfa), first auditorily, and then with accompanying orthography. Our results show 1) that participants only produced L1 [ɹ] when exposed to orthography, confirming that the use of Mandarin rhotic in L2 speech is orthographically driven; and 2) that even at the initial stage the substitution with Mandarin [ɹ] occurs al-

264

most exclusively in coda position, reminiscent of L2 learners [Zhou, 2017, Liu, 2018].

## Bibliography

Bene Bassetti. Orthography affects second language speech: Double letters and geminate production in English. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 43(11): 1835–1842, November 2017. ISSN 1939-1285, 0278-7393. DOI: 10.1037/xlm0000417. URL <http://doi.apa.org/getdoi.cfm?doi=10.1037/xlm0000417>.

Paul Boersma. Some listener-oriented accounts of h-aspiré in French. *Lingua*, 117(12):1989–2054, December 2007. ISSN 00243841. DOI: 10.1016/j.lingua.2006.11.004. URL <https://linkinghub.elsevier.com/retrieve/pii/S0024384106002191>.

Paola Escudero and Karin Wanrooij. The Effect of L1 Orthography on Non-native Vowel Perception. *Language and Speech*, 53(3): 343–365, September 2010. ISSN 0023-8309, 1756-6053. DOI: 10.1177/0023830910371447. URL <http://journals.sagepub.com/doi/10.1177/0023830910371447>.

Silke Hamann. Ghost phonemes in second languages – How orthography can create contrasts without perceptual correlates, 2018.

Silke Hamann and Ilaria E. Colombo. A formal account of the interaction of orthography and perception: English intervocalic consonants borrowed into Italian. *Natural Language & Linguistic Theory*, 35(3):683–714, August 2017. ISSN 0167-806X, 1573-0859. DOI: 10.1007/s11049-017-9362-3. URL <http://link.springer.com/10.1007/s11049-017-9362-3>.

Wen Liu. Aquisição da Vibrante Simples [r] pelos Alunos Chineses Aprendentes de Português como Língua Estrangeira, 2018.

Chao Zhou. Contributo para o estudo da aquisição das consoantes líquidas do português europeu por aprendentes chineses, 2017.

Chao Zhou and Silke Hamann. Cross-Linguistic Interaction Between Phonological Categorization and Orthography Predicts Prosodic Effects in the Acquisition of Portuguese Liquids by L1-Mandarin Learners. In *Interspeech 2020*, pages 4486–4490. ISCA, October 2020. DOI: 10.21437/Interspeech.2020-2689. URL <https://>

300 [www.isca-speech.org/archive/interspeech\\_2020/zhou20h\\_](http://www.isca-speech.org/archive/interspeech_2020/zhou20h_)  
301 [interspeech.html](http://www.isca-speech.org/archive/interspeech_2020/zhou20h_).



302 PART II:

303 ADVANCES IN SPEECH AND

304 L2 PROCESSING

Fairy tales are more than true: not because  
they tell us that dragons exist, but because  
they tell us dragons can be beaten.

C.K. CHESTERTON

# An Analysis of the Development of the Rhythm of English-L2 by Brazilian Learners through Rhythmic Metrics and Acoustic Parameters

LEONARDO ANTONIO SILVA TEIXEIRA<sup>1</sup>, RONALDO MANGUEIRA LIMA JR.<sup>1</sup>

1 . Universidade Federal do Ceará

## Abstract

The aim of this study is to describe and discuss the development of L2 English rhythm by Brazilian learners through rhythmic metrics and prosodic-acoustic parameters that characterize the oral production of these learners at different stages of L2 development. Five Brazilian learners of English-L2 were recorded reading a text in English at the beginning of their college studies in English Language Teaching, and again four semesters later, after having taken two English phonology courses. They were also recorded reading a version of the text translated into Portuguese. Besides the learners, five native speakers of North American English were recorded reading the same text in English. Data were manually segmented into vowel units (V), consonant (C), vowel-vowel (VV), sentences (S) and higher prosodic units - chunks (CH) in PRAAT [Boersma and Weenink, 2019], and the parameters were automatically by means of a script. Data were statistically treated via R [?]through the implementation of mixed-effects regression models. Results placed Brazilian Portuguese and English-L1 in different rhythmic spaces, as predicted by the literature; in the durational dimension, the metrics positioned the English-L2 of the first recording far from both English-L1 and Brazilian Portuguese; in the f0 and intensity dimensions, however, the acoustic parameters placed the English-L2 of the first

recording closer to Brazilian Portuguese. In both dimensions, the English-L2 of the subsequent recording was closer to English-L1, suggesting a developmental route towards the target language. The results also suggest positive effects of the explicit teaching of pronunciation.

## Introduction

Regarding research in non-native language (L2) development, there seems to be greater emphasis on segmental aspects rather than prosodic ones [Li and Post, 2014, Thomson and Derwing, 2015]. This tendency is also reflected in L2 acquisition models [Flege et al., 2021, Best and Tyler, 2007], which emphasize segmental aspects, providing little support to the understanding of L2 prosodic development. There is also evidence that rhythm can influence the communication process in a global way, affecting degrees of perceived foreign accent and intelligibility [Silva Junior and Barbosa, 2020].

The scarcity of studies on the acquisition of L2 rhythm may be related to the difficulty of establishing the physical reality of such construct. There are at least three trends on research regarding linguistic rhythm. Lloyd James (1940), as cited in Abercrombie [1971], relied on the dichotomy Morse code versus machine gun to illustrate the perceptual difference between English and Spanish, respectively. Pike [1945] formalized that difference by proposing a rhythmic approach based on the type of units that stood up in such languages: stress-timed languages, for which interstress intervals would be the most prominent units, and syllable-timed languages, for which syllables would be such units. Later, Abercrombie [1971] proposed that those rhythmic units would be isochronous, that is, of the same duration. However, the isochrony paradigm proved to be empirically unsustainable since intervals of the same duration are not found in the acoustic signal [Cumming, 2010].

From the mid-90s, a second trend of studies in linguistic rhythm emerges, in which the rhythmic patterns are investigated by means of the durational characteristics of the reference intervals (vowels, consonants, syllables, etc.), which can be computed by statistical indexes called rhythmic metrics. Ramus et al. [1999] proposed the standard deviation of the duration of consonantal intervals ( $\Delta C$ ) and the percentual of the total duration of the utterance composed of vowel intervals (%V). Those metrics were able to spatially discriminate languages considered syllable-timed (French, Spanish, Italian and Catalan), stress-timed (English, Polish and Dutch) and mora-timed languages

(Japanese) [Ladefoged, 1975] on a plane with  $\Delta C$  and %V on each axis, as can be seen in Figure 1, reproduced from the original paper:

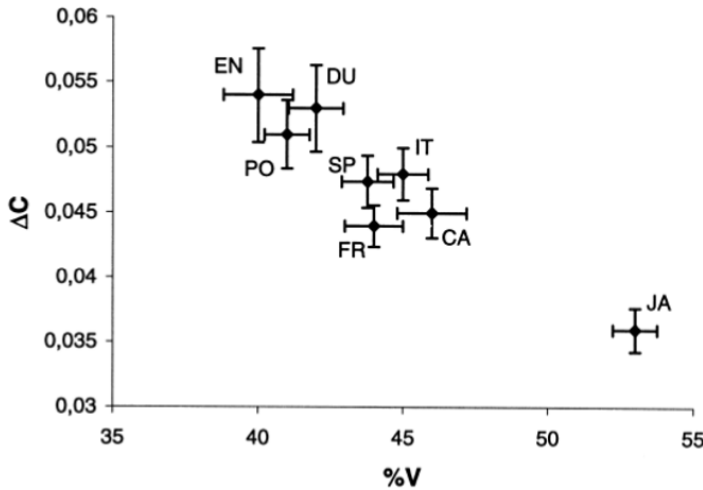


Figure 1: Distribution of languages over the (%V,  $\Delta C$ ) plane. Error bars represent 1 standard error. Source: Ramus et al. [1999, p. 273].

The present study follows a third trend on research on linguistic rhythm, which define it as a function of the distribution of prominent elements in the acoustic signal, which involves several acoustic dimensions – duration, fundamental frequency ( $f_0$ ) and intensity, and may be influenced by the native language of the speaker [Cumming, 2010, Fuchs, 2016, Silva Junior and Barbosa, 2020]. Thus, this study was guided by the following questions: (i) how do the metrics and acoustic parameters place North American English-L1, Brazilian English-L2, and the Brazilian Portuguese (BP)-L1 in the rhythmic space? (ii) What is the influence of the rhythm of BP-L1 on the development of English-L2 of learners? (iii) What is the effect of explicit pronunciation teaching on learners' English-L2 rhythm development? The following hypotheses were raised: (i) PB-L1, English-L1 and English-L2 are rhythmically different systems; (ii) there will be rhythmic differences between the English-L2 of the speakers in the two different stages of development whose recordings were analyzed; (iii) the English-L2 of the first recording should be more dissimilar to English-L1 due to L1 transfer and lack of explicit instruction.

## Methods

As for the participants, the experimental group was composed of five BP-L1 speakers, who were also learners of English-L2. They were all

college students of English Language Teaching, being four men and one woman, aged between 18 to 24. The control group comprised five English-L1 speakers, all Canadians, being one man and four women, aged 23-34. Four corpora of oral production were analyzed in this study: English-L1, PB-L1, English-L2 (1) and English-L2 (4). The data of English-L2 were obtained by means of recordings of the Brazilian learners reading the first paragraph of a text in two different moments, before and after completing courses in English Phonetics and Phonology. Those were the first and fourth recording made so they are referred to as English-L2 (1) and English-L2 (4). The data of English-L1 resulted from the reading of the same text by the control group. Finally, the Portuguese-L1 data came from the reading of the Portuguese version of the text by the Brazilian learners. The recordings took place in a silent room with a cardioid Shure MX150B lapel microphone connected to a Zoom 4HnSP recorder. The audio was captured in mono, with a sampling rate of 44.1 kHz, and saved in wav format.

Data were manually segmented into vowel units (V), consonant (C), vowel-vowel (VV), that is, the interval between the acoustic onset of a vowel and the onset of the adjacent one, sentences (S) and higher prosodic units - chunks (CH) in PRAAT [?], and the script Metrics & Acoustics Extractor [Silva Junior and Barbosa, 2020] was used to extract the parameters. Following Silva Junior and Barbosa [2020], the term metric(s) is used in this research to refer to the duration-based parameters, and the term acoustic parameter(s) refers to the f0, speech rate and intensive-related ones. The table below presents a summary of the metrics and acoustic parameters analyzed in this study and the types of segments they compute:

Data were then statistically treated via R [?]through the implementation of mixed-effects regression models, adopting language and semester as predictor variables, and rhythmic metrics and acoustic parameters as response variables.

## Results

In this section, we present some of the significant results, based on the mixed-effects regression models adjusted for each metric and acoustic parameter, and the boxplots and bidimensional planes, in which the effect of each corpus (independent variables) on the significant metrics and acoustic measures (the dependent variables) can be visually inspected.

<sup>1</sup> See Fuchs [2016] for a comprehensive account of rhythmic metrics.

<sup>2</sup> Standard deviation of the segment duration divided by the mean, multiplied by 100.

<sup>3</sup> Mean of the differences between successive segments.

<sup>4</sup> Mean of the differences between successive segments divided by their sum, multiplied by 100.

<sup>5</sup> Mean of pairwise quotients of adjacent segment durations, where the duration of the shorter is divided by the duration of the longer one and multiplied by 100.

<sup>6</sup> Mean of the differences between successive segments where the duration of each segment is normalised through division by the mean of all segments' durations.

<sup>7</sup> Mean of the differences between successive segments where the durations are normalised by z-transformation.

METRICS <sup>1</sup>		ACOUSTIC PAREMETERS	
Parameters	Segment of application	Parameters	Segment of application
Percentual (%)	V, C	f0 median	S, CH
Standard-deviation ( $\Delta$ )	V,C (V ou C), VV	f0 peak	S, CH
Variation coefficient (Varco) <sup>2</sup>	V,C (V ou C), VV	f0 minimum	S, CH
Raw pairwise variability index (r-PVI) <sup>3</sup>	V,C (V ou C), VV	f0 standard deviation	S, CH
Normalized pairwise variability index (n-PVI) <sup>4</sup>	V,C (V ou C), VV	f0 skewness	S, CH
Rhythm ratio (RR) <sup>5</sup>	V,C (V ou C), VV	Mean of f0 first derivative ( $\mu\Delta 1$ -f0)	S, CH
Variability index (VI) <sup>6</sup>	V,C (V ou C), VV	Standard deviation of f0 first derivative ( $\sigma\Delta 1$ -f0)	S, CH
Yet another rhythm determination (z-score duration) (YARD) <sup>7</sup>	V,C (V ou C), VV	Skewness of f0 first derivative ( $sk\Delta 1$ -f0)	
		Speech rate (SR)	VV, S, CH
		f0 rate (f0-R)	S, CH
		Spectral emphasis	S, CH
		Mean of normalized syllable-peak duration ( $\mu dur$ -Sil)	VV, S, CH
		Mean duration of pauses ( $\mu dur$ -#)	S, CH

Table 1: Rhythm metrics and prosodic-acoustic parameters analyzed in this study

## Metrics

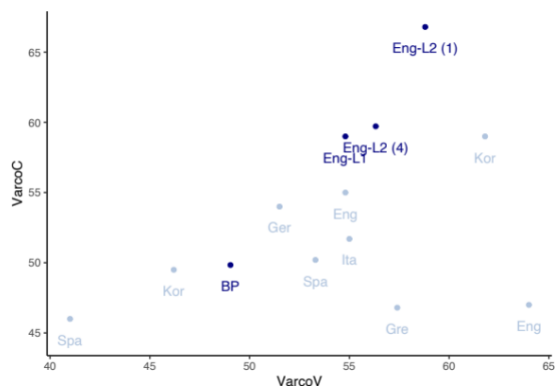
Twenty out of the thirty employed metrics reached statistical significance for at least two of the (inter) languages. One example of the mixed-regression models that were implemented via R can be seen in table 2, which were adjusted for the standard deviation of the duration of consonantal intervals ( $\Delta C$ ) and the percentual of vocalic intervals (%V).

Predictors	$\Delta C$	%V	p			
	Estimates	CI				
(Intercept)	46.48	36.17 – 56.79	<0.001	48.78	46.02 – 51.55	<0.001
Lang [Eng-L1]	21.94	7.35 – 36.52	0.004	-9.66	-12.97 – -6.35	<0.001
Lang [Eng-L2 (1)]	58.71	44.13 – 73.30	<0.001	-11.92	-14.23 – -9.61	<0.001
Lang [Eng-L2 (4)]	37.61	23.02 – 52.19	<0.001	-9.28	-11.47 – -7.09	<0.001

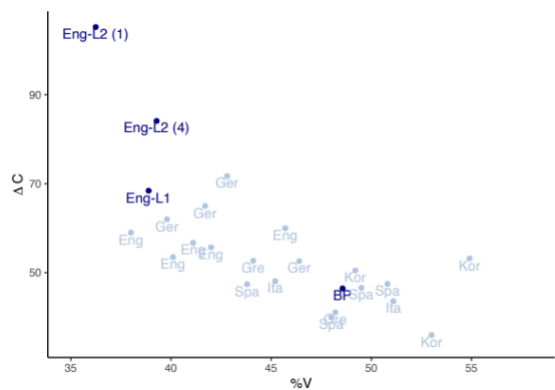
Figure 2 shows the distribution of the 4 corpora over the planes formed by the pairs  $\Delta C$ -%V and VarcoC-VarcoV in comparison to the data reviewed and obtained by Arvaniti (2012).

As can be seen in figure 2 (a), English-L1 presented greater standard deviation of consonantal intervals duration ( $\Delta C_{Eng-L1} = 68.41$ ) compared to BP ( $\Delta C_{BP} = 46.48$ ); and BP presented greater proportion of the utterance composed of vowel intervals ( $\%V_{BP} = 48.56$ ) compared to English-L1 ( $\%V_{Eng-L1} = 38.88$ ). The data of English-L2 (1) were positioned far from the two native languages, scoring  $\Delta C$  values quite high ( $\Delta C_{Eng-L2(1)} = 105.19$ ) and the lowest proportion of vowel segments ( $\%V_{Eng-L2(1)} = 36.24$ ). On the other hand, English-L2 (4) values were much closer to English-L1 in relation to both axis ( $\Delta C_{Eng-L2(4)} = 84.08$ ;  $\%V_{Eng-L2(4)} = 39.28$ ). As for VarcoC-VarcoV (Figure 2), English-L1, English-L2 (1), English-L2 (4) and BP were distributed analogously to the plane  $\Delta C$ -%V, with the BP data recording the lowest values for both the VarcoC (49.84) and VarcoV (49.04) axis, and English-L2 (1) presenting the highest scores both in the VarcoC axis (66.8) and in relation to the VarcoV axis (58.8). The fact that English-L2(1) assumed values far from the L1 (BP) indicates no objective transference of durational prosodic patterns to the learners' interlanguages. On the other hand, the approximation between English-L2(4) and English-L1 indicates a possible effect of explicit instruction, among other factors, that may have influenced the temporal (re)organization of the learners' speech towards the prosodic patterns of the target language.

Table 2: Coefficients, confidence intervals (95%) and p-Values for the two linear mixed-effect regression models adjusted for  $\Delta C$  and %V. models:  $\Delta C$  Lang + (1|Chunk) + (1|Speaker) and percV Lang + (1|Chunk) + (1 / Speaker).



(a)



(b)

Figure 2: Present study data (dark blue) amid all the data reviewed and obtained by Arvaniti (2012) (light blue) for  $\Delta C$  - %V (Figure 2 (a) and VarcoC-VarcoV (Figure 2 (b)), in which Eng = English, Ger = German, Gre = Greek, Spa = Spanish, UI = Italian, Kor = Korean. Source: Teixeira and Lima Jr. (2021).



Regarding the data from Arvaniti (2012), BP grouped with languages considered more syllable-timed, that is, with more durational regularity among the segments of reference, such as Spanish and Italian. English-L1 results were also consistent with the literature, gathering with the results for English and German from other studies, which are considered languages with more stress-timing tendency.

The hierarchy of values for VarcoV-VarcoC and  $\Delta C$ - $\%V$  illustrates the dominant positioning pattern for the significant metrics, as can be seen in table 3: ([ + stress-timed] English-L2 (1) > English-L2 (4) > English-L1 > BP [ + syllable-timed]), except for  $\%V$  and RR, whose higher values indicate a tendency towards syllable-timing.

Metric	BP	English-L1	English-L2(1)	English-L2(4)
$\%V$	48.56 (3.16)	38.88 (4.96)	36.24 (5.36)	46.48(8.02)
$\%C$	51.44 (3,16)	61.12 (4.96)	63.76 (5.36)	68.416(14.55)
$\Delta V$	40.08 (10.81)	41.16 (11.82)	51.81 (12.79)	105.192(32.6)
$\Delta C$	46.48 (8.02)	68.41 (14.55)	105.192 (32.6)	84.088(36.51)
$\Delta S^*$	133.4 (45.27)	198.53 (77.44)	217.46(97.65)	184.75 (56.72)
VarcoV	49.04 (8.49)	54.80 (12.52)	58.80 (11.30)	56.32(14.81)
VarcoC	49.84 (7,00)	59 (11.89)	66.80 (13.69)	59.72(22.02)
rPVI-V	65.1 (11.71)	70.74 (14.84)	96.98 (19.27)	81.21(15.75)
rPVI-C	48.22 (8.91)	86.16 (20.23)	116.84 (46.64)	88.7(21.69)
rPVI-VC	64.73 (18.72)	83.58 (11.12)	114.4 (38.18)	89.1(14.53)
rPVI-S	102.85 (40.24)	130.66 (33.59)	176.27 (90.03)	137.96(44.62)
nPVI-C	53.96 (7.21)	68.56 (11.72)	72.36 (12.14)	64.84(11.46)
nPVI-VC	59.76 (7.69)	68.96 (11.09)	72.6 (9.40)	65.08(8.12)
RR-C	61.17 (4.36)	53.13 (6.29)	50.97 (6.59)	54.59(6.42)
RR-VC	58.07 (4.35)	52.8 (5.65)	50.91 (4.97)	54.55(4.59)
VI-V	0.818 (0.166)	0.981 (0.322)	1.128 (0.302)	0.894(0.188)
VI-V	0.830 (0.037)	0.924 (0.049)	0.929 (0.052)	0.934(0.059)
VI-VC	0.684 (0.101)	0.834 (0.157)	0.859 (0.120)	0.746(0.116)
VI-S	0.516 (0.120)	0.606 (0.136)	0.615 (0.160)	0.538(0.126)
YARD-VC	0.717 (0.150)	0.695 (0.133)	0.869 (0.113)	0.848(0.123)

Table 3: Absolute means for the statistically significant metrics and standard deviation (between parentheses) for BP, English-L1, English-L2 (1), English-L2(1) and English-L2(4).  
 \* S stands for the phonetic syllable, which is the vowel-vowel (VV) unit.

#### Acoustic Parameters

Five out of the twelve employed acoustic parameters reached statistical significance:  $f0_{peak}$ ,  $\sigma f0$ ,  $\sigma \Delta f0$ , spectral emphasis (emph) and speech rate (SR).

As for the standard deviation of  $f_0$  (Figure 3.1), English-L1 presented the highest standard deviation among the corpora analyzed ( $\sigma f_0 \text{Eng-L1} = 3.79$ ), followed by English-L2(4) ( $\sigma f_0 \text{Eng-L2(4)} = 3.34$ ), English-L2(1) ( $\sigma f_0 \text{Eng-L2(1)} = 2.71$ ) and BP ( $\sigma f_0 \text{BP} = 2.62$ ). The results for this parameter suggest a gradual prosodic development of the learners towards the  $f_0$  variation patterns of the target language. The standard deviation of  $f_0$  first derivative ( $\sigma \Delta 1-f_0$ ) (Figure 3.2) was also successful in the separation of the L1s and captured a similar course of development to that found by  $\sigma f_0$ . The highest mean was scored by English-L1 ( $\sigma \Delta 1-f_0 \text{Eng-L1} = 5.51$ ), the lowest mean was scored by BP ( $\sigma \Delta 1-f_0 \text{BP} = 3.61$ ). The interlanguages registered intermediate values, but the mean English-L2(4) was much closer to English-L1 ( $\sigma \Delta 1-f_0 \text{Eng-L2(1)} = 3.73 < \sigma \Delta 1-f_0 \text{Eng-L2(4)} = 4.61$ ).

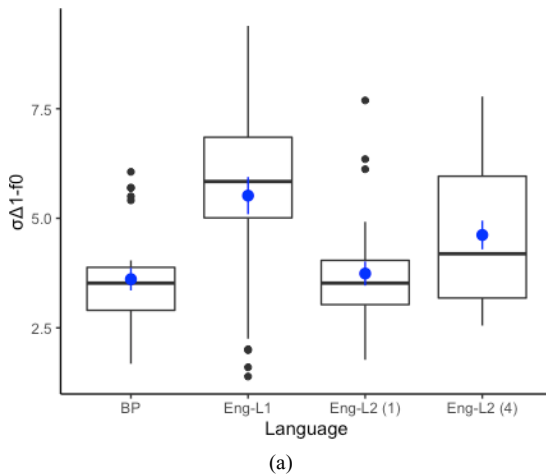
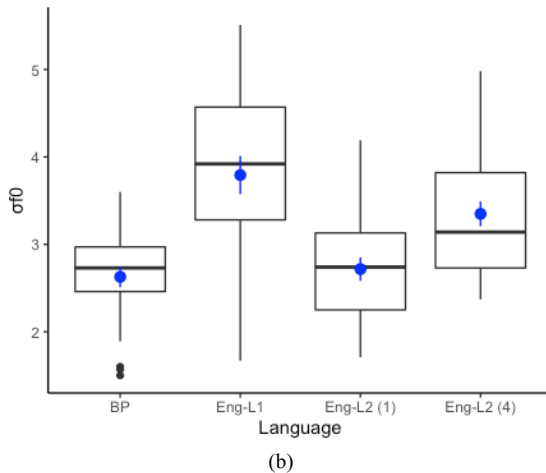


Figure 3: Boxplots of the means  $\sigma f_0$  (Figure 3.1) and  $\sigma \Delta 1-f_0$  (Figure 3.2) for English-L1, English-L2(1), English-L2(4) and BP. The blue dots and lines represent the means and standard errors respectively.



The results for the  $f_0$  dimension must be interpreted with caution, since there was an unbalance between male and female participants in both groups (control group: 1 male, 4 female; experimental group: 4 males, 1 female). In fact, the correlation between  $f_0$  and sex is evident when individual results are taken into consideration. For instance, in the experimental group, it was observed that participant N, the only female, is the one that had the highest  $f_0$  peak (97.16), as well as the widest scopes of  $f_0$  ( $f_0$  peak minus  $f_0$  min) for BP (17.18) and English-L2(4) (19.06). There was also a smaller variation between the male learners  $f_0$  scope of English-L2(1) ( $A = 12.28$ ;  $F = 15.68$ ;  $K = 15.5$ ;  $L = 13.37$ ) and English-L2(4) ( $A = 12.81$ ;  $F = 15.09$ ;  $K = 14.43$ ;  $L = 15.1$ ), in comparison to the variation of the female participant, who went from 15.59 to 19.06 in the last recording.

In the dimension of intensity, as visually demonstrated in Figure 4.1, spectral emphasis was able to separate the L1s, with the highest mean for English-L1 among the analyzed corpora ( $\text{emphEng-L1} = 4.34$ ), which was higher than PB ( $\text{emphPB} = 2.73$ ). If we consider works that show the correlation between spectral emphasis and phrasal stress [?], this result suggests that native English speakers make more effort as an acoustic clue in stress marking than Portuguese speakers. Regarding the interlanguages, English-L2 (1) obtained the lowest mean of spectral emphasis, very close to BP values, ( $\text{emphEng-L2(1)} = 2.56$ ), and English-L2 (4) got much closer to English-L1 ( $\text{emphEng-L2 (4)} = 3.23$ ). This indicates L1 transfer at the intensity dimension, and a tendency towards the prosodic patterns of English-L1 in the last recording.

As expected, the L1s presented higher speech rates, with BP registering a higher mean compared to English-L1 ( $\text{SRPB} = 5.22 > \text{SREng-L1} = 4.43$ ). In addition, English-L2 (1) presented the lowest speech rate among the corpora analyzed ( $\text{SREng-L2(1)} = 3.59$ ) and English-L2(4) registered a slightly higher mean, closer to English-L1 ( $\text{SREng-L2(4)} = 3.74$ ). The increase in the speech rate of the interlanguages between the first and last recording may be related to the effects of explicit instruction.

## Discussion

The metrics and parameters positioned BP, English-L1, English-L2 (1) and English-L2 (4) as rhythmically different systems. There were differences between the English-L2 of the speakers in the two different stages of development and different developmental paths were captured as function of the dimension of prominence. This developmental path

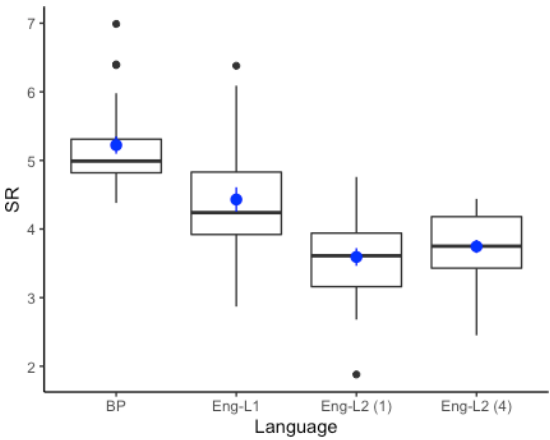
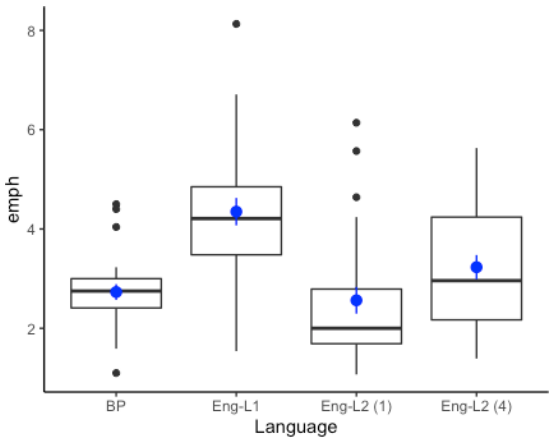


Figure 4: Boxplots of the means spectral emphasis (Figure 4 (a)) and speech rate (Figure 4 (b)) for English-L1, English-L2(1), English-L2(4) and BP. The blue dots and lines represent the means and standard errors respectively.



(b)

is consistent with the definition of interlanguage that presents itself as a relatively independent system of L1 and L2 [Li and Post, 2014], and with the non-linearity of the L2 development process [Lima Jr and Alves, 2019].

At the durational level, the dominant distribution pattern ([+ stress-timed] English-L2 (1) > English-L2 (4) > English-L1 > BP [+ syllable-timed]), with English-L2(1) assuming the highest means among the four corpora, and English-L2(4) getting closer values to English-L1. One possible explanation for such behavior for English-L2(1) is that learners may have mobilized a process of dissimilation of phonetic categories, displaying exaggerated durational values to maintain the distinction between L1 and L2, similarly to what is predicted by the Speech Learning Model for the segmental level [?Flege et al., 2021].

At the f0 dimension, the dominant distribution pattern ([+ f0 variability] English-L1 > English-L2 (4) > English-L2 (1) > BP [- f0 variability]) placed English-L1 with the highest means among the 4 corpora and BP with the lowest means, which suggests English native speakers mobilize more complex and varied f0 contours in speech. L1 transfer was more salient at the f0 level and seems to be more persistent among men, which adds to Urbani [2012]. A tendency towards the f0 prosodic patterns of the target language was also identified at this level and could be an effect of explicit instruction. This effect may have also influenced the greater speech rate ([+ speech rate] BP > English-L1 > English-L2 (4) > English-L2 (1) [- speech rate]) and spectral emphasis ([+ spectral emph] English-L1 > English-L2 (4) > BP > English-L2 (1) [- spectral emph]) for English-L2 (4), suggesting more fluency of the learners, and an overall improvement in marking syllable stress, respectively.

## Conclusion

The metrics and acoustic parameters confirmed the first hypothesis that North American English-L1, the Brazilian English-L2, and the BP-L1 are rhythmically different systems. We confirmed the hypothesis that the data of English-L2 (1) would be more dissimilar in relation to English-L1 compared to the data of English-L2 (4), but orthogonal patterns of rhythmic development seem to coexist as a function of the different dimensions of prominence. Nevertheless, the approximation between the means of English-L2(4) and English-L1 in all dimensions suggest positive effects of explicit pronunciation in the development of prosodic features by learners of non-native languages. As future work,

we intend to expand the analyzed corpora including the recordings of the 2nd and 3rd semesters as well as the other paragraphs of the text, and to analyze the correlation between the metrics and acoustic parameters and perceived degrees of foreign accent, intelligibility, and comprehensibility.

## Acknowledgment

This study integrates a project that is partially financed by CNPq, process 438823/2018-4.

## Bibliography

David Abercrombie. Elements of general phonetics. Aldine Atherton, Chicago, 1971. URL [https://archive.org/details/elementsofgenera0000aber\\_u2o1](https://archive.org/details/elementsofgenera0000aber_u2o1). OCLC: 1256469368.

Catherine T. Best and Michael D. Tyler. Nonnative and second-language speech perception: Commonalities and complementarities. In Ocke-Schwen Bohn and Murray J. Munro, editors, *Language Learning & Language Teaching*, volume 17, pages 13–34. John Benjamins Publishing Company, Amsterdam, 2007. ISBN 9789027219732 9789027292872. DOI: 10.1075/llt.17.07bes. URL <https://benjamins.com/catalog/llt.17.07bes>.

Paul Boersma and David Weenink. Praat: doing phonetics by computer. 6.0.20, 2019. URL <http://www.praat.org>.

Ruth Elizabeth Cumming. Speech rhythm: the language-specific integration of pitch and duration. PhD thesis, University of Cambridge, Cambridge, 11 2010.

James Emil Flege, Katsura Aoyama, and Ocke-Schwen Bohn. The revised speech learning model (SLM-r) applied. In Raa-tree Wayland, editor, *Second Language Speech Learning*, pages 84–118. Cambridge University Press, 1 edition, February 2021. ISBN 9781108886901 9781108840637 9781108814614. DOI: 10.1017/9781108886901.003. URL [https://www.cambridge.org/core/product/identifier/9781108886901%23CN-bp-2/type/book\\_part](https://www.cambridge.org/core/product/identifier/9781108886901%23CN-bp-2/type/book_part).

Robert Fuchs. Speech rhythm in varieties of English: evidence from educated Indian English and British English. PhD thesis, 2016. OCLC: 1084748564.

Peter Ladefoged. Course in Phonetics. Houghton Mifflin Harcourt P,  
New York, 1975. ISBN 9780155151802.

Aike Li and Brechtje Post. L2 acquisition of prosodic properties of  
speech rhythm: evidence from L1 mandarin and german learn-  
ers of english. *Studies in Second Language Acquisition*, 36  
(2):223–255, June 2014. ISSN 0272-2631, 1470-1545. DOI:  
10.1017/S0272263113000752. URL [https://www.cambridge.org/core/product/identifier/S0272263113000752/type/journal\\_article](https://www.cambridge.org/core/product/identifier/S0272263113000752/type/journal_article).

Ronaldo Manguera Lima Jr and Ubiratã Kickhöfel Alves. A dynamic  
perspective on L2 pronunciation development: bridging research  
and communicative teaching practice. *Revista do GEL*, 16(2):  
27–56, December 2019. ISSN 1984-591X, 1806-4906. DOI:  
10.21165/gel.v16i2.2417. URL <https://revistas.gel.org.br/rg/article/view/2417>.

Kenneth L. Pike. The intonation of American English. University of  
Michigan publications Linguistics. University of Michigan Publica-  
tions, Ann Arbor, 1945.

Franck Ramus, Marina Nespor, and Jacques Mehler. Correlates of  
linguistic rhythm in the speech signal. *Cognition*, 73(3):265–  
292, December 1999. ISSN 00100277. DOI: 10.1016/S0010-  
0277(99)00058-X. URL <https://linkinghub.elsevier.com/retrieve/pii/S001002779900058X>.

Leônidas José da Silva Junior and Plínio Almeida Barbosa. Speech  
rhythm of english as L2: an investigation of prosodic variables  
on the production of Brazilian Portuguese speakers. *Journal of  
Speech Sciences*, 8(2):37–57, August 2020. ISSN 2236-9740.  
DOI: 10.20396/joss.v8i2.14996. URL <https://econtents.bc.unicamp.br/inpec/index.php/joss/article/view/14996>.

R. I. Thomson and T. M. Derwing. The effectiveness of L2 pro-  
nunciation instruction: a narrative review. *Applied Linguistics*,  
36(3):326–344, July 2015. ISSN 0142-6001, 1477-450X. DOI:  
10.1093/applin/amu076. URL <https://academic.oup.com/applij/article-lookup/doi/10.1093/applin/amu076>.

Martina Urbani. Pitch range in L1/L2 English. An analysis of F0  
using LTD and linguistic measures. *Methodological Perspectives on  
Second Language Prosody*, pages 79–83, 2012.

# Change-Point Analysis in language development: a study of voice onset time production in a mul- tilingual system

LAURA CASTILHOS SCHERESCHEWSKY<sup>1,2</sup>, UBIRATÃ KICKHÖFEL  
ALVES<sup>1,3</sup>

1 . Universidade Federal do Rio Grande do Sul

2 . CAPES

3 . CNPq

## Introduction

According to Complex Dynamic Systems Theory (CDST) <sup>1</sup>, when it comes to multilingual development, we need to think about the interconnectedness of the system components. Departing from this assumption, we follow Kupske's concept of language attrition<sup>2</sup>, which characterizes this phenomenon as the force resulting from the contact of two bodies, in this case, two languages, that are in constant movement [Kupske, 2016, p. 39–40]. This concept embraces the CDST premise that change is inherent to development. Thus, if the system is in constant movement, we may find it in continuous change in a given state, and language variability is expected to be found. Sometimes, the system may go through significant changes that exceed its current state [van Dijk and van Geert, 2007]. If these particular changes lead to the reorganization of the system as a whole (as in the emergence of a new attractor state), we call them 'phase transitions' or 'phase shifts'. According to Hepford [2020], this new attractor state is not necessarily something new to learners, as it "could be a language form that they

<sup>1</sup> ; ; ; and

<sup>2</sup> In this paper, we do not differentiate 'language attrition' from 'language transfer' or 'language drift'. We will use the three terms interchangeably.



are exposed to regularly but have not had the cognitive ability to adapt to, or an event that pushes a learner to adapt and self-organize resulting in using a new form” [Hepford, 2020, p. 162–163].

Based on the aforementioned assumptions, this study aims to investigate the phenomena that occur in the development of the additional languages of trilingual speakers, native speakers of Brazilian Portuguese (BP-L1) and non-native speakers of English (L2) and French (L3). Specifically, this longitudinal study analyzes, over a period of three months (with 12 weekly datapoints), the development of the production of Voice Onset Time (VOT), observing possible phase shifts through change-point analyses [Taylor, 2000, cf.] provided by the Change-point analyzer v.2.3 software [TAYLOR ENTERPRISES]. The study included a period of pedagogical intervention to accelerate the development of the positive VOT pattern with the characteristic aspiration of English. This teaching intervention took place over six explicit pronunciation instruction sessions, conducted in the weeks of datapoints 4 to 9. We aimed to discuss to what extent the accelerated development of an L2 with a typologically different VOT pattern causes changes in the development of the L3 and L1 subsystems, as well as show the inter-relation of the two additional languages over time.

According to the literature [Schereschewsky, 2021, cf.], from the study of VOT, we can observe the multidirectionality of transfer and the adaptability and the self-organization of language subsystems. Therefore, this study intends to provide empirical and theoretical input into a larger understanding of these aspects. This may shed light on the development of additional languages in the light of CDST. As this is essentially a theory about change, we aim to raise issues such as language development and its ongoing ”process” in time [Lowie and Verspoor, 2015, 2019, cf.], the interconnectivity of typologically different subsystems, data variability, and the emergence of new attractor states and phase shifts.

## Method

As addressed in the previous section, the main goal of this study was to inferentially verify possible phase shifts in VOT patterns, in each of the language subsystems, especially after the beginning of explicit pronunciation instruction in English. For that, we carried out change-point analyses [Taylor, 2000, cf.].

In order to achieve our goal, we proposed a methodology in which changes and interactions among the subsystems of multilingual speak-

ers could be investigated through accelerating L2 VOT development. The experiment was built with a longitudinal design in an A-B-A format [Hiver and Al-Hoorie, 2020, cf.], with 12 datapoints, which were intersected in the midpoints with 6 sessions of explicit pronunciation instruction in English. This intervention took place between the weeks referring to datapoints 4 and 9, and all instructional sessions were conducted with a communicative approach.

In this study, we replicated different process-oriented analyses to encompass and address variability [de Bot et al., 2013], conducting the same experiment with five participants from different backgrounds, different ages, with different proficiency levels in their additional languages, and different routines. Due to space restrictions, in this paper we will focus on the results from one particular participant<sup>3</sup>, who was 24 years old at the time, a graduate student who worked as a French teacher. This participant took a self-evaluation test [?] and graded herself a 6 in English and a 10 in French<sup>4</sup>.

The participant was presented with three different reading tasks. In each data collection session, she received 23 carrier sentences (repeated three times each) with 18 target words with /p/, /t/ and /k/ in word-initial position and 5 distractor words. The BP and English instruments were the same as in Kupske (2016), for both matters of consistency and comparisons of results with previous studies. We also used the same methodological control as Kupske in the development of the French instrument [Schereschewsky, 2021, cf.]. Because she received the same target words, the order of the carrier sentences was randomized and the distractor words were changed each week. This study was conducted during the pandemic of COVID-19, so the participant accomplished the experiment in an individual setting and was asked to complete each task taking time intervals between them.

All audio recordings from the reading tasks were analyzed acoustically in the Praat v.6.1.16 software [?]. Due to time and space restrictions, only the absolute values of VOT production were considered. As for VOT measurements, similar criteria to previous works were used: selecting the voiceless interval between the burst of the stop consonant and the first regular pulse of the following vowel.

As for the statistical analyses, the participants' developmental trajectories were plotted, considering minimum, maximum, and mean values from the tokens of each stop consonant in each datapoint. Following that, change-point analyses [Taylor, 2000, Steenbeek et al., 2012, Baba and Nitta, 2014, Han and Hiver, 2018, Englhardt et al., 2020, Henry et al., 2021, cf.] were conducted. Change-point analysis is an infer-

<sup>3</sup> This participant is referred to as Participant 5 in previous works from this project. For more information on the other participants, see Schereschewsky [2021].

<sup>4</sup> It is interesting to note that her L3 was more active than her L2, because she worked as a French teacher, even though she had started studying English before she even started learning French.

748 ential method that uses resampling and cumulative sums to identify a  
749 pattern shift, or the point of change, in a set of longitudinal data.

750 The Change-Point Analyzer software is able to detect several lon-  
751 gitudinal changes. By running a fast analysis of cumulative sums and  
752 bootstrapping, for each change in pattern, the software provides prac-  
753 tical information, including the confidence level, which indicates a  
754 probability that a change has actually occurred, and the confidence in-  
755 terval, which indicates when that change has occurred. Figure 1 shows  
756 the first output tab from the software, with the change-point analysis  
757 visual plot.

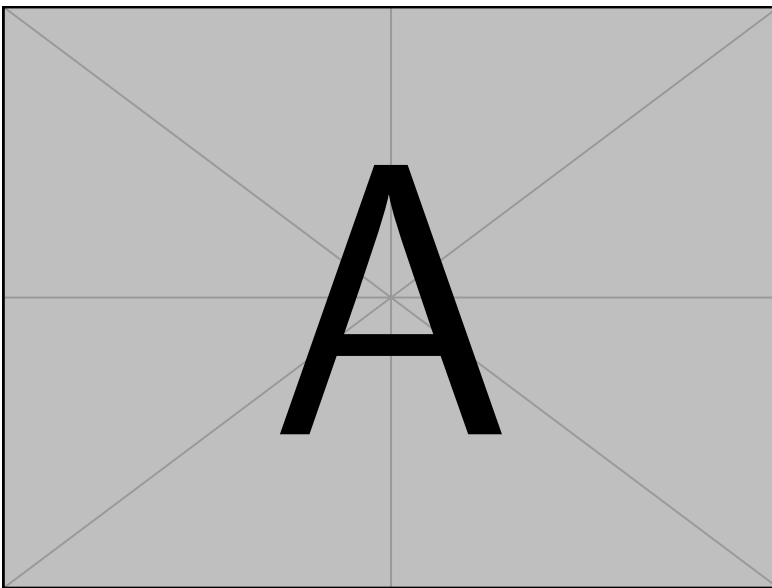


Figure 5: Outputs -  
change-point analysis  
visual plot.

758 In Figure 5, the bold black line in the graph represents the raw data  
759 of the mean VOT values of [p] in Participant 5's L1 over the 12 col-  
760 lection points. The dark blue lines represent the amplitude range of  
761 the control limits, that is, the maximum range of variation in which  
762 the values can fluctuate, assuming that no change has occurred (if the  
763 black line exceeds the control limits, we will have a first indicative that  
764 a change has taken place, which may simply be an outlier or an indica-  
765 tion of an actual phase shift). The lighter blue background represents  
766 the area that should contain all values varying within the control limits.  
767 The displacement of this area in light blue at the bottom of the graph  
768 actually indicates a phase shift, as the average values within the first

segment show a sudden change, starting to vary in a different range, represented by the second segment of the area in lighter blue. Figure 6 shows the second output tab from the software, with the significant changes of the set of data.

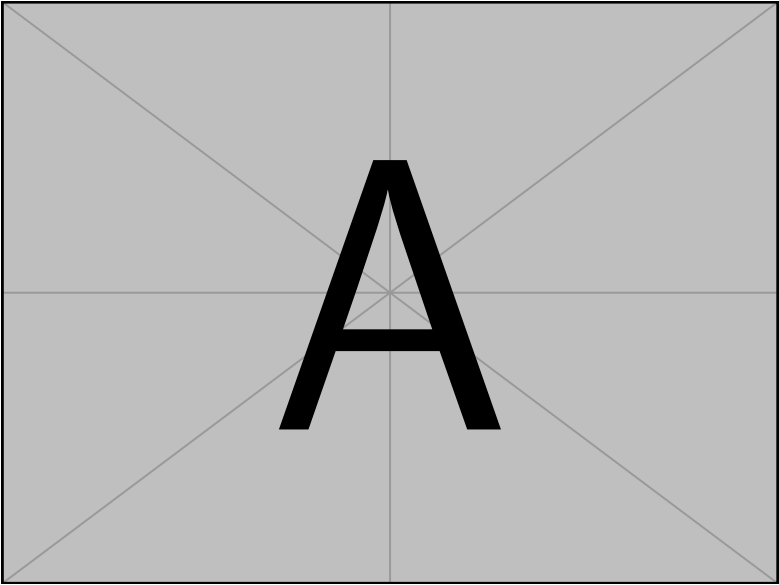


Figure 6: Outputs - table of significant changes.

The table in Figure 6 indicates the estimated point of change to another phase, in this case, in Datapoint #5, with a confidence level of 97%<sup>5</sup>, indicated by the confidence interval (which, in this case, points exactly to session 5). Next, the table indicates the values before and after the change, that is, the average values of variation in the first phase (considering the average of all inputs within this first phase), which go from 34.25ms to 46.29ms in the second phase. Finally, the level of change indicates its importance. In this particular example, the Level 1 change indicates that this was the first significant change identified by the software in the first analysis run of the data. Other change levels may appear, depending on how many phase changes are identified and whether these are significant. Finally, Figure 7 shows the third output tab from the software, with the visual chart showing the cumulative sums.

The chart in Figure 7 represents the cumulative sum analyses (CUSUM). According to Taylor [2000, p. 6], "they are the cumulative sums of differences between the values and the average". These dif-

<sup>5</sup> The Change-point Analyzer only presents, in the outputs, intervals that have at least 95% confidence. The more spaced the confidence interval, the lower the confidence level for a change to have occurred at the point identified by the software.

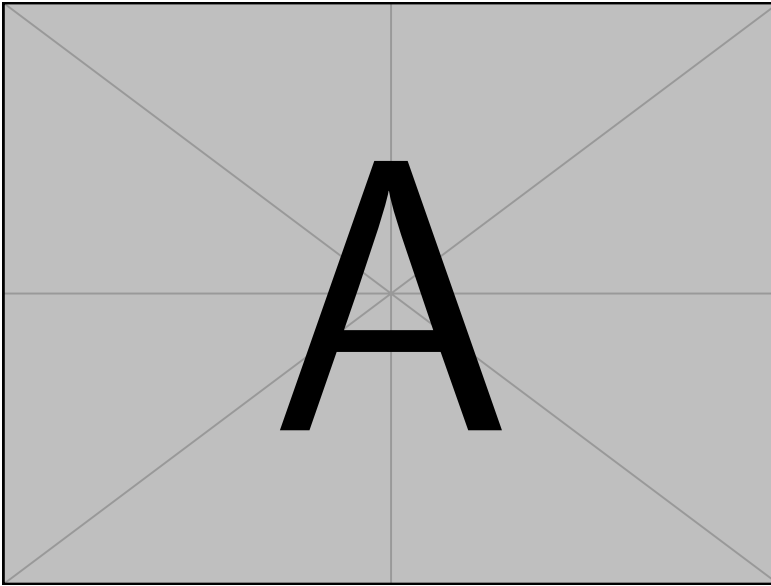


Figure 7: Outputs - table of significant changes.

790 ferences sum to zero so the cumulative sum always ends at zero. Thus,  
791 in the CUSUM graph from this data, a downward sloping line can be  
792 seen, indicating that the values in that period have a tendency to be  
793 below the general average, until there is a change in the direction of the  
794 line, starting to have an upward slope, indicating that values from that  
795 portion of the graph tend to be above the overall mean. We should also  
796 look at the shaded background of the chart, which indicates whether  
797 and where there has been a significant change in the slope of the line,  
798 referring to the table with confidence intervals. Another important in-  
799 formation brought by the CUSUM chart is that the straighter the line,  
800 regardless of its direction (up or down), the greater the certainty that no  
801 change occurred in that period. On the other hand, the more curved the  
802 line is, as is the case between points 4 and 9, the greater the possibility  
803 that other changes (from other levels) have taken place.

## 804 Results

805 As previously mentioned, change-point analyses are used to identify  
806 the points at which a pattern change occurs in a longitudinal dataset.  
807 Thus, change-point analyses help identify the developmental stages of  
808 VOT production, checking attractor states in phases of relative stability

in each language. In this section, only a summary of the significant changes and the most relevant charts for the discussion will be presented. A table will be presented for each language, with the significant results split by consonant (/p, t, k/), of the data collection session the change took place, the confidence interval, the confidence level (in percentages), the mean values before and after the change (which is related to the averages of variation of values within the control limits of each phase) and the level of change (degree of importance in the analysis by the software). Table 1 shows the results of Brazilian Portuguese-L1.

Stop	Measure	Session	Conf.level	From	To	Shift
[p]	Mean	5	97%	34,252	46,289	↗
[p]	Max	5	95%	69,097	84,571	↗
[t]	Mean	11	100%	40,181	32,775	↘
[t]	Max	11	99%	76,253	48,585	↘
[k]	Mean	4	96%	59,073	75,531	↗

Table 4: Change-point analysis of BP-L1.

First, we emphasize the interconnectedness of the language subsystems, which makes it possible for a native language to change, even if it is typologically different from a language that underwent an intervention, as we found significant phase changes in the production of VOT in Portuguese-L1 in the three stops.

For [p], we found a Level 1 phase shift in the means in Datapoint 5, when the averages change from 34.25ms to 46.29ms, and a Level 3 change in maximums around Datapoint 5, when the averages increase from 69.1ms to 84.57ms.

For [t], in which we also found significant phase changes in the averages and maximum instances, the data are somewhat more interesting. In both measures, Level 2 phase changes are found, in which the VOT decreases in duration. For the means, phases shifted from an average of 40.18ms to 32.78ms. For the maximum instances, they shifted from 76.25ms to 48.59ms. However, as both changes are Level 2 and this new phase with shorter VOT measures only starts by the end of the analyzed period, around Datapoint 11, it is also necessary to visually analyze these data, since there is also the possibility that another (non-significant) change occurred in previous datapoints. Figure 8 shows the change-point analysis plots for the two measures.

The graphs of the means and the maximum instances of [t] show that the two measures presented a very similar behavior during the analyzed period. Comparing the initial and final points of the mea-

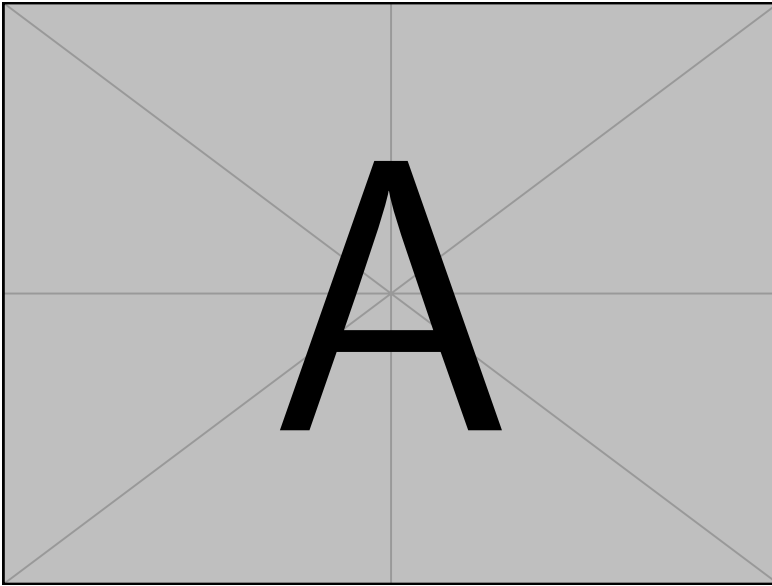


Figure 8: Change-point analysis of means and maximum instances of [t] in Portuguese-L1.

842 surements, there is a clear trend towards a decrease in the descriptive  
843 values of VOT, which is in accordance with the phase shift found with  
844 a decrease of averages. However, there is also a very clear indication  
845 that the data may have undergone another phase shift, around Datapoint  
846 3, where the VOT appears to have increased in duration. The CUSUMs  
847 graph shows the possibility of another phase shift due to the sudden  
848 change in the direction of the cumulative sums line on the third data-  
849 point. However, the software did not verify this change as significant  
850 to include it in the outputs. What was included in the outputs was a  
851 significant Level 1 phase shift in the means of [k], where there was an  
852 increase in the averages from 59.07ms to 75.53ms in Datapoint 4, thus  
853 after the beginning of the intervention, once again showing that even  
854 the subsystem of a typologically different language is subject to change  
855 as a result of another one changing.

856 The English-L2 results also bring valuable data to the discussion.  
857 For [p], there is a significant change in Datapoint 4, the first after the  
858 start of the intervention. When it comes to the means, the Level 2  
859 phase shift occurs when the average changes from 54.55ms to 95.42ms.  
860 For the maximums, the Level 1 change occurs with an increase of the  
861 averages from a phase of 104.37ms to 142.81, with very high values of  
862 VOT production for a bilabial stop.

Stop	Measure	Session	Conf.level	From	To	Shift
[p]	Mean	4	97%	54,553	95,416	↗
[p]	Max	4	94%	104,37	142,81	↗
[t]	Min	12	91%	40,409	26,24	↘
[t]	Mean	4	94%	62,407	105,89	↗
[t]	Mean	11	93%	105,89	80,07	↘
[t]	Max	4	100%	110,1	147,2	↗
[k]	Mean	4	99%	82,73	112,96	↗
[k]	Max	5	91%	130,65	157,73	↗

Table 5: Change-point analysis of English-L2.

For [t], we found significant phase changes for the three analyzed measures, but each measure presented a different result. For minimums, for instance, a Level 3 phase shift occurs around Datapoint 12, with a decrease in averages from 40.41ms to 26.24ms. With such a large confidence interval, which covers the entire intervention until the end of the study, in addition to the fact that it is a Level 3 change, there remains a possibility of another, less significant phase shift, in some other datapoint in that interval. For the means, two significant phase changes were identified, one of Level 1, in Datapoint 4, with an increase in the averages from 62.40ms to 105.89ms, and one of Level 2, at the end of the study, around Datapoint 11, this time with a decrease in averages (much like what happens in her L1) from 105.89ms to 80.07ms, a higher average than in the initial phase. The maximums of [t] present a third pattern of behavior, with a Level 3 phase shift around Datapoint 4, with an increase from 110.1ms to 147.2ms in the later phase.

Figure 9 shows the plots of the change-point analyses of the three analyzed measures of [t] in English-L2. Although the three measures show completely different behaviors in the outputs, as evidenced by the first graph of each one, the CUSUM graphs of the three are very similar, indicating sudden changes in the slope of the cumulatives sums line at least twice on each<sup>6</sup>. This pattern would be indicative of at least three distinct phases during the study, in which the first phase change would represent an increase in the average VOT values, and the second a slight decrease, as we verified in the means of [t], almost always involving the same datapoints. For minimums and maximums, however, a possible second phase change was not significant, leaving the observation only for a qualitative discussion.

Finally, for [k], we found significant phase shifts in the means and in the maximums, both indicating an increase in VOT values. For the

<sup>6</sup> A downward-sloping CUSUM line indicates values below the overall average, while an upward-sloping line indicates values above the overall average. A change in the slope of the line represents a change in trend and, being a significant change, corresponds to a phase shift.



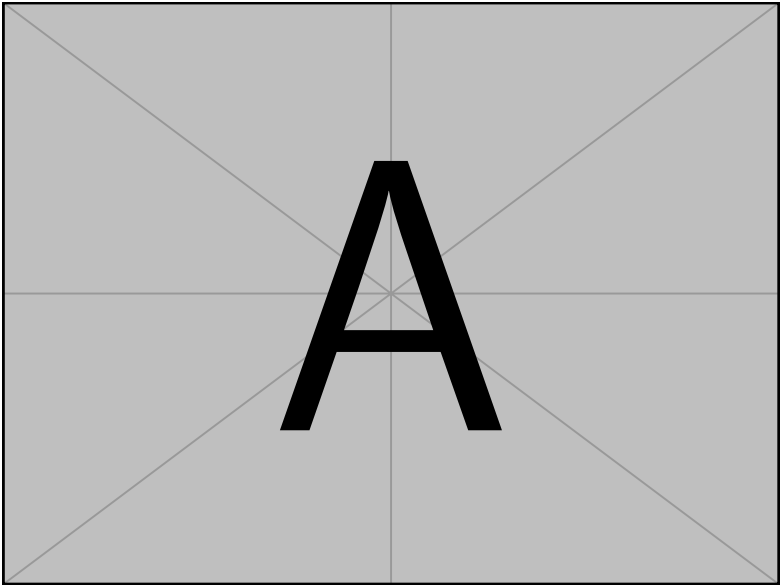


Figure 9: Change-point analyses of the minimums, means and maximums of [t] in English-L2.

means, the change occurred in Datapoint 4, where a new phase went from an average of 82.73ms to 112.96ms. For the maximums, the change occurred around Datapoint 5, when the results show an increase of the averages from 130.65ms to 157.73ms. Overall, all these English-L2 data are extremely valuable in showing the influence of explicit instruction in the development of new attractor states, that is, new phases developing a non-native positive VOT pattern with long-lag aspiration. Furthermore, these data highlight change as an inherent characteristic of a developing system, showing that the language remains in motion even after the end of an intervention.

Stop	Measure	Session	Conf.level	From	To	Shift
[p]	Mean	4	97%	30,987	40,782	↗
[p]	Max	4	96%	60,697	78,477	↗
[t]	Mean	4	99%	33,437	40,15	↗
[t]	Max	4	95%	58,91	70,327	↗
[k]	Mean	6	92%	59,452	66,256	↗

Table 6: Change-point analysis of French-L3.

Once again, we can observe significant phase shifts in the three consonants of a language subsystem that is typologically different from the language that received explicit instruction during the intervention,

showing the interconnectivity of the system as a whole. Interestingly, all the identified changes present new phases with an increase in the VOT values in the French language. For [p], the means and maximums undergo phase shifts in Datapoint 4. For the means, the Level 2 shift showed a change in the averages from 30.99ms to 40.78ms. For the maximums, the Level 1 shift showed changed averages from 60.68ms to 78.48ms. For visualization purposes, the graphs referring to the phase shifts in the maximums of [p] are in the Figure 10.

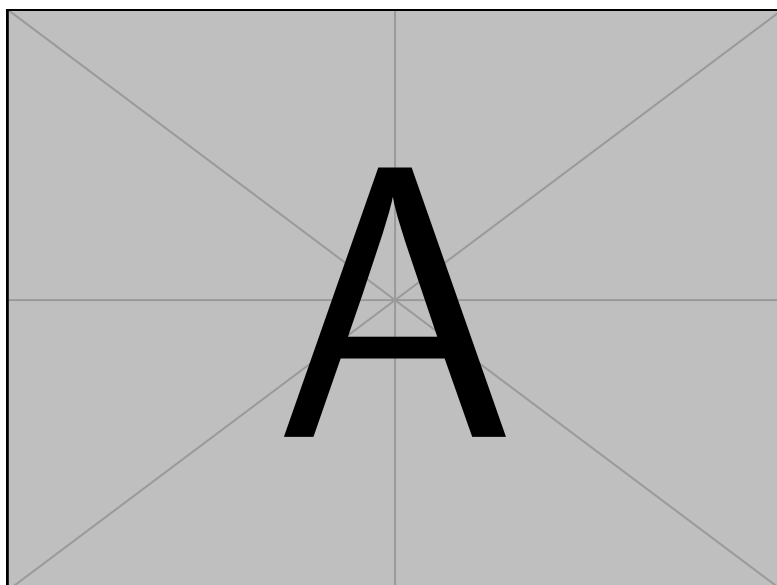


Figure 10: Change-point analyses of the maximums of [p] in French-L3.

For [t], Level 1 phase shifts were also identified in the means and maximums, occurring around Datapoint 4. For the means, the phase shift indicates an increase in average from 33.44ms to 40.15ms, and in the maximums, from 58.91ms to 70.33ms.

On the other hand, finally, we only found a significant phase shift in the means of [k]. The Level 1 change was identified around Datapoint 6, indicating an increase in average from 59.45ms to 66.26ms. Again, we reiterate the subsystem's ability to change under the influence of changes in other subsystems, especially the L2, since the new phases were always identified from the beginning of the instruction period in that language.

## Final Considerations

We highlight the relevance of change-point analyses in verifying the emergence of new developmental phases and we emphasize that change-point analyses allow us to identify more than one change in each language subsystem, as shown in the L2 data. Considering that changes in a multilingual system are constant and that even attractor states are not permanent, an analysis of this sort provides valuable information on the process of multilingual development. As shown in our results, languages are entangled and interconnected in a multilingual system, and they influence one another. Finally, we hope to contribute to the area of language development in the light of Complex Dynamic Systems Theory. Discussing methods of analysis that verify developmental changes is always necessary. Specifically, change-point analyses help to identify the emergence of new stages of development. As a non-linear process, we acknowledge the fluctuations of the VOT values in the new developmental phase, as it probably refers to a less strong attractor state, and even the emergence of a third phase, different from both the initial phase and the phase under the influence of the pedagogical intervention. These results show that language and learning are constantly changing, demonstrating the relevance of an approach via CDST, given that this, after all, constitutes a theory essentially about change.

## Bibliography

- Kyoko Baba and Ryo Nitta. Phase Transitions in Development of Writing Fluency From a Complex Dynamic Systems Perspective: Phase Transition in Development of Writing Fluency. *Language Learning*, 64(1):1–35, March 2014. ISSN 00238333. DOI: 10.1111/lang.12033. URL <https://onlinelibrary.wiley.com/doi/10.1111/lang.12033>.
- Kees de Bot, Wander Lowie, Steven L. Thorne, and Marjolijn Verhoeven. Dynamic systems theory as a comprehensive theory of second language development. In María del Pilar García Mayo, María Junkal Gutierrez Mangado, and María Martínez Adrián, editors, *AILA Applied Linguistics Series*, volume 9, pages 199–220. John Benjamins Publishing Company, Amsterdam, 2013. ISBN 9789027205254 9789027205285 9789027272225. DOI: 10.1075/aals.9.13ch10. URL <https://benjamins.com/catalog/aals.9.13ch10>.

- 963 Adrian Englhardt, Jens Willkomm, Martin Schäler, and Klemens  
964 Böhm. Improving semantic change analysis by combining word  
965 embeddings and word frequencies. *International Journal on*  
966 *Digital Libraries*, 21(3):247–264, September 2020. ISSN 1432-  
967 5012, 1432-1300. DOI: 10.1007/s00799-019-00271-6. URL  
968 <http://link.springer.com/10.1007/s00799-019-00271-6>.
- 969 Jiwon Han and Phil Hiver. Genre-based L2 writing instruction  
970 and writing-specific psychological factors: The dynamics of  
971 change. *Journal of Second Language Writing*, 40:44–59, June  
972 2018. ISSN 10603743. DOI: 10.1016/j.jslw.2018.03.001.  
973 URL [https://linkinghub.elsevier.com/retrieve/pii/](https://linkinghub.elsevier.com/retrieve/pii/S1060374317304642)  
974 [S1060374317304642](https://linkinghub.elsevier.com/retrieve/pii/S1060374317304642).
- 975 Alastair Henry, Cecilia Thorsen, and Peter D. MacIntyre. Willing-  
976 ness to communicate in a multilingual context: part one, a time-  
977 serial study of developmental dynamics. *Journal of Multilingual*  
978 *and Multicultural Development*, pages 1–20, June 2021. ISSN  
979 0143-4632, 1747-7557. DOI: 10.1080/01434632.2021.1931248.  
980 URL [https://www.tandfonline.com/doi/full/10.1080/](https://www.tandfonline.com/doi/full/10.1080/01434632.2021.1931248)  
981 [01434632.2021.1931248](https://www.tandfonline.com/doi/full/10.1080/01434632.2021.1931248).
- 982 Elizabeth Hepford. Chapter 7. The elusive phase shift: Capturing  
983 changes in L2 writing development and interaction between the  
984 cognitive and social ecosystems. In Gary G. Fogal and Mar-  
985 jolijn H. Verspoor, editors, *Language Learning & Language*  
986 *Teaching*, volume 54, pages 161–182. John Benjamins Publish-  
987 ing Company, Amsterdam, June 2020. ISBN 9789027205575  
988 9789027205582 9789027261144. DOI: 10.1075/llt.54.07hep.  
989 URL <https://benjamins.com/catalog/llt.54.07hep>.
- 990 Phil Hiver and Ali H. Al-Hoorie. Research methods for complexity  
991 theory in applied linguistics. Number 137 in *Second language*  
992 *acquisition*. Multilingual Matters, Bristol ; Blue Ridge Summit,  
993 2020. ISBN 9781788925730 9781788925747.
- 994 Felipe Flores Kupske. Imigração, atrito e complexidade: a pro-  
995 dução das oclusivas surdas iniciais do inglês e do português por  
996 sul-brasileiros residentes em Londres. PhD thesis, Universidade  
997 Federal do Rio Grande do Sul, Porto Alegre, 2016.
- 998 Diane Larsen-Freeman and Lynne Cameron. *Complex systems and*  
999 *applied linguistics*. Oxford applied linguistics. Oxford university  
1000 press, Oxford, 2008. ISBN 9780194422444.

- 1001 Wander Lowie and Marjolijn Verspoor. Variability and variation  
1002 in second language acquisition orders: A dynamic reevaluation:  
1003 Variability in acquisition orders: Dst. *Language Learning*, 65(1):  
1004 63–88, March 2015. ISSN 00238333. DOI: 10.1111/lang.12093.  
1005 URL [https://onlinelibrary.wiley.com/doi/10.1111/lang.](https://onlinelibrary.wiley.com/doi/10.1111/lang.12093)  
1006 12093.
- 1007 Wander M. Lowie and Marjolijn H. Verspoor. Individual Differences  
1008 and the Ergodicity Problem: Individual Differences and Ergodicity.  
1009 *Language Learning*, 69:184–206, March 2019. ISSN 00238333.  
1010 DOI: 10.1111/lang.12324. URL [https://onlinelibrary.](https://onlinelibrary.wiley.com/doi/10.1111/lang.12324)  
1011 [wiley.com/doi/10.1111/lang.12324](https://onlinelibrary.wiley.com/doi/10.1111/lang.12324).
- 1012 Laura Castilhos Schereschewsky. Desenvolvimento de voice onset time  
1013 em sistemas multilíngues (português - 11, inglês - 12 e francês - 13):  
1014 discussões dinâmicas a partir de diferentes metodologias de análise  
1015 de processo. Master’s thesis, Universidade Federal do Rio Grande do  
1016 Sul, Porto Alegre, 2021.
- 1017 Henderien Steenbeek, Louise Jansen, and Paul van Geert. Scaffold-  
1018 ing dynamics and the emergence of problematic learning trajec-  
1019 tories. *Learning and Individual Differences*, 22(1):64–75, Febru-  
1020 ary 2012. ISSN 10416080. DOI: 10.1016/j.lindif.2011.11.014.  
1021 URL [https://linkinghub.elsevier.com/retrieve/pii/](https://linkinghub.elsevier.com/retrieve/pii/S1041608011001646)  
1022 [S1041608011001646](https://linkinghub.elsevier.com/retrieve/pii/S1041608011001646).
- 1023 Wayne Taylor. Change-Point Analysis: A Powerful New Tool For  
1024 Detecting Changes, April 2000. URL [https://variation.com/](https://variation.com/change-point-analysis-a-powerful-new-tool-for-detecting-changes/)  
1025 [change-point-analysis-a-powerful-new-tool-for-detecting-changes/](https://variation.com/change-point-analysis-a-powerful-new-tool-for-detecting-changes/).
- 1026 TAYLOR ENTERPRISES. Change-Point Analyzer. URL [https:](https://variation.com/product/change-point-analyzer/)  
1027 [//variation.com/product/change-point-analyzer/](https://variation.com/product/change-point-analyzer/).
- 1028 Marijn van Dijk and Paul van Geert. Wobbles, humps and sudden  
1029 jumps: a case study of continuity, discontinuity and variability  
1030 in early language development. *Infant and Child Development*,  
1031 16(1):7–33, February 2007. ISSN 15227227, 15227219. DOI:  
1032 10.1002/icd.506. URL [https://onlinelibrary.wiley.com/](https://onlinelibrary.wiley.com/doi/10.1002/icd.506)  
1033 [doi/10.1002/icd.506](https://onlinelibrary.wiley.com/doi/10.1002/icd.506).

# 1034 Production of English [Cs] clusters by Brazilian 1035 speakers: effects of orthography, phonological 1036 environment and task type

1037 WELLINGTON ARAUJO MENDES JUNIOR<sup>1</sup>

1038 1 . Federal University of Minas Gerais

1039 This study examines the effects of orthography, phonological en-  
1040 vironment and task type in the production of English [Cs] clusters  
1041 by Brazilian Portuguese (BP) speakers. Two orthographic patterns  
1042 were examined for English nouns whose plural is pronounced as a  
1043 (stop + sibilant) cluster. One of the patterns presents two consonants  
1044 word-finally - cups, cats, ducks - whereas the other one presents a  
1045 silent vowel <e> between two consonants: grapes, plates, cakes.  
1046 The goal was to assess whether these different orthographic patterns  
1047 would trigger the production of an epenthetic vowel. Additionally,  
1048 it was assessed whether different phonological environments would  
1049 influence the voicing property of the final sibilant. As it is known,  
1050 word-final English sibilants are prone to progressive assimilation (e.g.  
1051 cups [kʌps], bags [bægz]), rather than regressive assimilation – as it  
1052 occurs in BP (mês [mes], mês anterior [mez ã.te.ri.'or]). An experi-  
1053 ment was designed to test the production of [Cs] clusters in English  
1054 nouns and in BP forms undergoing sound change. Harmonics-to-noise  
1055 ratio (HNR) was used to measure sibilant voicing, whereas the pres-  
1056 ence of epenthetic vowels was assessed categorically. Results showed  
1057 that English learners are more likely to pronounce a vowel when the  
1058 orthographic pattern is <Ces> rather than <Cs>, and this occurs  
1059 regardless of the visual presentation of the words. Moreover, HNR  
1060 rates showed that fully voiced sibilants tend to occur in L2 English

when the consonant is both preceded and followed by a vowel. These findings are discussed in light of the Exemplar Model in L2 Phonology (EML2P) [Cristófar-Silva and Guimaraes, 2021, Mendes Jr. and Cristófar-Silva, 2022 in press]. The analysis based on the EML2P showed that robust patterns from the L1 are adopted in L2, including fine phonetic detail that reflects subphonemic properties.

## Introduction

Traditional phonological models assume that English plural suffixes and third person singular present forms are subject to a phonological rule. The underlying representation for regular plural and 3<sup>rd</sup> person singular present is assumed to be /z/ [Hayes, 2011]. A progressive assimilation rule predicts that if a vowel or a voiced consonant precedes /z/, the output is [z], as in *dogs* [dɒgz], *trees* [triːz] and *pies* [paɪz]. If a voiceless consonant precedes /z/, it surfaces as [s], as in *cups* [kʌps], *cats* [kæts] and *ducks* [dʌks]. Finally, if an alveolar fricative or an affricate precedes the sibilant, the outcome is [ɪz], as in *inbuses* [bʌsɪz], *quizzes* [kwɪzɪz] and *watches* [wɒtʃɪz]. However, when we consider the orthography of English plural forms, two possible spellings are associated with the aforementioned sound patterns. Nouns can either end in a consonant followed by the letter <s>, as in *books* and *jobs* or the by letters <es>, as in *cakes* and *cubes*. Brazilian Portuguese, on the other hand, presents mainly the <Ces> pattern, as it occurs in *cheques* and *clubes*, whereas only some few nouns present the <Cs> pattern: *biceps*, *forceps*, *volts*.

As a consequence of an ongoing sound change, word-final [Cs] clusters<sup>7</sup> are currently very productive in some Brazilian Portuguese plural forms: *crepes*, *potes*, *cheques* [Soares, 2016]. The alternation between [Cs] ~ [Cis] word-finally in BP follows from the reduction and eventual loss of unstressed high front vowels when flanked between a consonant and a word-final sibilant. It seems that such alternation also applies to plural forms produced by Brazilian speakers of L2 English, as in *incakes* [keɪks] ~ [ˈkeɪ.kis].

This paper intends to investigate [Cs] clusters in English regular plural forms (e.g. *cups* [kʌps], *grapes* [greɪps]) produced by Brazilian speakers of L2 English in an attempt to address the question of whether an ongoing sound change from the L1 plays a role in L2 learning. Additionally, we aim to assess how orthographic and phonological representations are related. Studies on the relationship between orthography and phonology have increased in recent years [Rafat, 2015, Hamann and Colombo, 2017, Zhou, 2021]. The main research questions in this topic

<sup>7</sup> For the purpose of the present discussion, we refer to [Cs] as any (consonant + sibilant) sequence. However, as it will be discussed later, the sibilant may be either voiced or voiceless.

aim to explain how L2 learners mediate the relationship between the already known phonological and orthographical knowledge from the L1 in order to build an L2. Thus, an important question we pose is whether different orthographic patterns trigger different pronunciations of [Cs] clusters in L2 English.

This paper is organized as follows. The next section reviews studies on the production of English [Cs] clusters by Brazilian speakers. The third section describes the methodology adopted in this study. The fourth section discusses our findings and is followed by the conclusions.

## Production of English [Cs] clusters by Brazilian speakers

Several works have addressed the relationship between orthography and the pronunciation of L2 English forms by Brazilian speakers. One of the main concerns have been to assess whether the presence of an epenthetic vowel in L2 English is influenced by a letter corresponding to a vowel. Delatorre [2006] investigated the production of English [Cs] clusters that occur in past and participle forms by Brazilian speakers of L2 English (e.g. moved and robbed). Two epenthetic vowels were attested: one epenthetic vowel breaks up the word-internal consonant cluster and the other one prevents word-final consonants, as in asked[ˈas.ke.dʒi] and saved[ˈseɪ.ve.dʒi]. Delatorre [2006] claimed that the orthographic input, which was present in a reading task, favored higher rates of an epenthetic vowel, as opposed to a free speech task, which did not present any orthographic stimulus. Therefore, she argued that the orthographic input favored the presence of epenthetic vowels in the pronunciation of L2 English by Brazilian speakers.

Although she did not focus on the production of [Cs] clusters, Silveira [2007] also investigated the production of word-final epenthesis in Brazilian speakers of L2 English. She compared words whose final letter was a consonant (e.g. mad[mæd]) to words whose final letter was a silent <e> (e.g. made[mɛɪd]). Her results showed that words ending in a silent <e> presented higher rates of epenthesis than words that ended in a consonantal letter. Akin to Delatorre [2006], the results of Silveira [2007] showed that a reading task favored higher rates of epenthetic vowels than a free speech task, indicating that orthographic input (and the task type) contributed to the production of an epenthetic vowel.

Another case of epenthetic vowels reported in the literature involves word-final consonant and sibilant sequences, [Cs], which typically appear in regular plural and 3<sup>rd</sup> person singular present forms in English. It is known that Brazilian speakers of L2 English tended to insert an



1141 epenthetic vowel between two word-final consonants, as it occurs, for  
1142 example, in cakes [keiks]~ [ˈkei.kis] [ʔ]. Interestingly, works that con-  
1143 sidered 3<sup>rd</sup> person singular present and regular plural forms in English  
1144 spoken by Brazilian speakers did not account for an epenthetic vowel.  
1145 They were rather concerned with voice agreement.

1146 Zanfra [2013] studied sibilant voicing in L2 English by Brazilian  
1147 speakers. Although her focus was not specifically on plural forms, her  
1148 results shed some light on the current discussion. The author tested  
1149 whether the BP voicing assimilation rule involving adjacent segments  
1150 in word boundaries would apply in L2 English learners' productions.  
1151 Her results showed that sibilants tended to be voiced when followed  
1152 by a voiced consonant (e.g. The house backyard is huge) or by a vowel  
1153 (e.g. The mouse I saw is white). Conversely, a sibilant was voiceless  
1154 when the following context was a pause (e.g. I won't go if he goes.) or  
1155 a voiceless consonant (e.g. These pancakes are great). Zanfra [2013]  
1156 suggested that Brazilian speakers of L2 English transfer the  
1157 BP regressive assimilation rule into their L2 English.

1158 Fragozo [2017] investigated the voicing of sibilants in English regu-  
1159 lar plural forms and 3<sup>rd</sup> person singular presented by Brazilian speakers  
1160 of L2 English. She assessed the extents to which a sibilant would be  
1161 voiced after a voiced consonant, as in dogs or clubs, which would re-  
1162 flect the acquisition of a progressive assimilation rule from English.  
1163 Fragozo [2017] also examined words in context to verify if the regres-  
1164 sive assimilation rule, which applies to BP, would be transferred to L2  
1165 English. She found that voiced sibilants tended to follow the regressive  
1166 assimilation rule from BP, whereas the English progressive assimi-  
1167 lation rule had a very low rate in her data (0.6%). She argues that the  
1168 low rates of voiced sibilants [z] in L2 English by Brazilian speakers  
1169 follows from the fact that these consonants are only partially voiced in  
1170 English. Data from her control group of native speakers presented 44%  
1171 of expected voiced sibilants. Thus, as sibilants are partially voiced in  
1172 English, they would not be accessible in L2 English.

1173 Zanfra [2013] and Fragozo [2017] both investigated voicing agree-  
1174 ment within a rule-based approach where there would be a competition  
1175 between a regressive assimilation rule from BP and a progressive  
1176 assimilation rule from English. A question that arises from this as-  
1177 sumption is whether a rule that is transferred from the L1 to the L2  
1178 could change as time goes by. Another issue which is polemic lies on  
1179 the role played by orthography, as in hou<se> orbu<s> [Zanfra,  
1180 2013]. Orthography cannot be modelled within a rule-based approach  
1181 as it is not part of Grammar. Furthermore, the rule-based approach  
1182 adopted by Zanfra [2013] and Fragozo [2017] neglected the role played

by an epenthetic vowel that may intervene between the two word-final consonants, as in cakes[keiks] ~ [ˈkeɪ.kis]. Additionally, they did not account for the gradience of sibilant voicing.

Unlike previous works which adopt rule-based approaches, this paper models L2 phonology within an Exemplar Model by considering representation robustness and the role of fine phonetic detail in shaping mental representations. Within this proposal, orthography is modelled as part of the linguistic knowledge of literate speakers and sound patterns display a great range of variability and gradience.

## Methodology

A set of 36 plural nouns ending in a sequence of (stop + sibilant) were considered in BP. These words present a single orthographic pattern: <Ces>, as in cheques [ʃɛks] ~ [ˈʃɛ.kis] ‘cheques’. For the L2 English case study, a set of 36 words were selected, where 15 words display the orthographic pattern <Ces>, as in grapes [greɪps], and the other 21 words display the orthographic pattern <Cs>, as in maps [mæps].

The experiment comprised two tasks. The first one consisted of a picture-counting task in which participants were asked to count and name the items shown in the pictures. Short carrier sentences that did not include orthographic stimuli of the target words were given. The second trial consisted of a reading task. Initially, participants were asked to read 72 BP sentences aloud. Alike the picture-counting task, BP nouns in the reading task were followed by either a vowel or a voiceless consonant. On the other hand, L2 English nouns were followed by either a vowel or a pause. The overall number of syllables was controlled for both languages: 4 in English and 12 in BP, considering the deletion of the [i] vowel. Sentence-level intonation and the morphological class of each word were also controlled.

A group of six Brazilians studying at the Federal Center for Technological Education of Minas Gerais, in the city of Araxá, participated in this study<sup>8</sup>. All participants were high school students who had been taking English classes as part of the school’s curriculum for about one year. The group consisted of 3 males and 3 females and their ages ranged from 15 to 17. All participants displayed either B1 or B2 proficiency levels (intermediate learners) of the Common European Framework of Reference for Languages.

Due to the recent COVID-19 pandemic, all interactions were performed remotely. Experiments were recorded with the Open Broadcaster Software Studio at 48 kHz sampling rate. The obtained recordings were converted into WAVEform audio format by the soft-

<sup>8</sup> This research has been approved by the ethics committee from the Universidade Federal de Minas Gerais, reference number: CAAE: 15116119.9.0000.5149.

were Adobe Premiere 2020, which was able to maintain the same sampling rate as the original files. The average time to complete the experiment was 45 minutes. A total of 648 tokens were collected for the L2 English study. For the BP study, 432 tokens were collected. Samples were edited and manually annotated using Praat TextGrids [?].

Besides assessing the presence or absence of a vowel between [Cs] clusters, this research also considered the voice quality of word-final sibilants. In BP, only voiceless sibilants occur word-finally, unless a vowel follows it, to which a voiced sibilant occurs. In English, voiced and voiceless sibilants occur word-finally. When a vowel follows the sibilant, the voice quality remains as it formerly was (rather than changing as it occurs in BP). We posited that word-final voiceless sibilants would be favored in L2 English, as it is the more robust pattern in L1. We also posited that a voiced sibilant occurs at higher rates in an intervocalic position: [Cis] followed by a word-initial vowel.

Voicing was measured under Harmonics-to-noise ratio. Each token was extracted to a separate sound object and a harmonicity object was created, from which the mean harmonicity was calculated, hereafter the HNR. The details of its calculation can be found in Boersma [1993]. Harmonicity would seem to be a good measurement of voicing since vocal cord vibration produces “a complex periodic wave” [Johnson, 1997, p. 63]. Based on the discussion from Praat’s manual, higher values of HNR should correspond with higher voicing rates.

## Results

Consider Figure 11, which shows the rates of [Cs] in regular plural forms in BP and L2 English.

The leftmost column shows that regular plural forms in BP, whose orthography is <Ces>, presented 62% of a consonant followed by a sibilant: [Cs]. That means that when a letter <e> appears in the orthography of BP plural forms, a vowel is manifested in 38% of the cases. The two rightmost columns report data from English spoken by Brazilian speakers. When the orthography in the plural form is <Ces>, a consonant followed by a sibilant [Cs] occurred in 83% of the cases, whereas in the cases where the orthography was <Cs>, a consonant followed by a sibilant occurred in 96% of the productions. This result shows that the pronunciation of [Cs] is more recurrent when the orthography is <Cs> than when the orthography is <Ces> in regular plural forms in English. In other words, a vowel will appear at higher rates when the orthographic pattern is <Ces> than when it is <Cs>. Thus, it is more likely that a plural form as tapes will

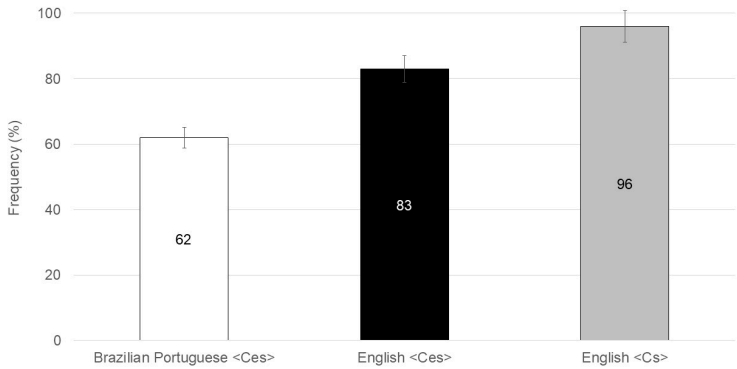


Figure 11: [Cs] rates by orthographic patterns.

have a vowel pronounced between the last two consonants than a plural form as maps. The difference between the data presented in the two rightmost columns is statistically significant for the orthographic patterns ( $\chi^2 = 36.113$ ,  $df = 1$ ,  $p < 0.01$ ). The explanation for such difference lies in the different orthographic patterns.

We also considered whether different tasks could favor the production of an epenthetic vowel. According to Delatorre [2006] and Silveira [2007], visual input favors such non-target productions. In our experiment, the picture-counting task had no orthographic visual input, whereas orthography was available in the reading task. If Delatorre [2006] and Silveira [2007] are correct, then we expect that vowels would occur at higher rates in the reading task than in the picture-counting task in our experiment. However, no statistically significant differences were found between the picture-counting task and the reading task ( $\chi^2 = 0.66$ ,  $df = 1$ ,  $p\text{-value} = 0.41$ ). This shows that it is the orthographic pattern rather the type of task that favors a vowel to occur in L2 English. Our claim is that once speakers are literate, orthography is part of their grammar, i.e., it has a permanent impact on mental representations. The EML2P model adopted in the current paper differs from Delatorre [2006], Silveira [2007] and Zangra [2013] rule-based approach mainly by assuming that orthography is part of linguistic knowledge and not external to it.

Another research question we posited regarded the voice quality of the word-final sibilant in [Cs] and [Cis]. This was the main issue considered by Zangra [2013], Fragozo [2017] within a rule-based approach. Their analysis claimed that voicing in L2 English did not achieve the

target rates due to constraints of BP distribution of sibilants and regressive assimilation. BP only presents voiceless sibilants word-finally. However, across word-boundaries, BP sibilants are voiced when followed by a voiced consonant or a vowel: *mês* [mes] ‘month’, *mês bonito* [mez 'bo.ni.tu] ‘beautiful month’, *mês anterior* [mez ɔ̃.te.ri.'or] ‘previous month’. In this paper, we offer an alternative view to the preceding rule-based approaches. Within the scope of the EML2P, it is suggested that generalizations from an ongoing sound change in BP phonology are transferred into L2 English, where phonetic detail plays an important role in shaping mental representations. Consider Figure 12.

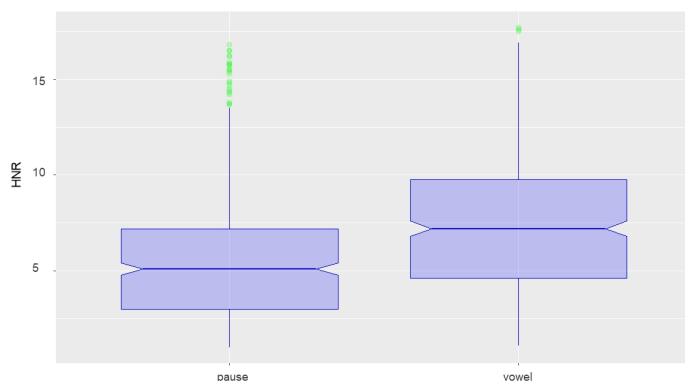


Figure 12: HNR per following phonetic environment in L2 English.

The boxplots in Figure 12 show harmonics-to-noise ratio per following phonetic environment in L2 English. We can see that when the sibilant is followed by a pause, it tends to be unvoiced, with HNR rates at around 5 decibels. Conversely, when the sibilant is followed by a vowel, voicing rates are higher. T-test results show that there is a significant difference in HNR between both following phonetic environments ( $t = -8.8153$ ,  $df = 821.37$ ,  $p\text{-value} < 0,01$ ). However, even though such environments seem to influence voicing rates of the final sibilant, these rates are still lower when compared to English target forms. To put it another way, nouns that should be pronounced with a word-final voiced sibilant present more unexpected voiceless sibilants than voiced ones. This can be accounted by the fact that only voiceless sibilants occur word-finally in BP. Learners are likely unaware of the fact the [z] should be voiced in accordance with the voice property of the preceding segment. Now consider Figure 13.

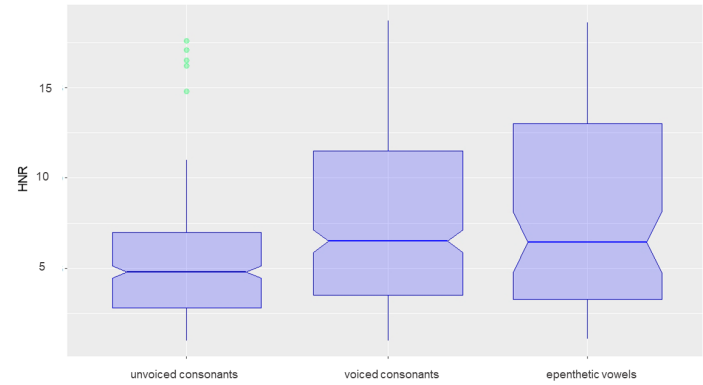


Figure 13: HNR per preceding phonetic environment in L2 English.

The boxplots in Figure 13 show harmonics-to-noise ratio per preceding phonetic environment in L2 English. Data is comprised of sibilants preceded either by an unvoiced consonant, a voiced consonant or an epenthetic vowel. At first sight, we can see the HNR rates are somewhat lower when the sibilant is preceded by an unvoiced consonant, and higher rates occur when the sibilant is followed by voiced consonants and epenthetic vowels. An analysis of variance (ANOVA) on these scores yielded significant variation among conditions:  $F(2,858)=59.99$ ,  $p < 0.001$ . A post-hoc Tukey test showed that the group comprised of unvoiced consonants differed significantly at  $p < 0.05$ ; the voiced consonants group was not significantly different from the epenthetic vowels group. This result suggests that epenthetic vowels contribute to higher rates of voicing as much as other voiced segments in L2 English. Finally, an interaction between both preceding and following phonetic environments was attested [ $F(2,858)=7.797$ ,  $p < 0.001$ ].

Our results throw some light on the line of research carried out by Zafra [2013] and Fragozo [2017], who investigated the sibilant voicing followed by a vowel within rule-based approaches. We account for the fact that low HNR (which reflect voiceless sibilants) is recurrent in regular plural forms in L2 English, as [s] is the most robust exemplar in word-final position in BP. We also account for the fact that the pattern [Cis] favors a voiced sibilant in L2 English, as voiced sibilants are favored in similar contexts in BP (i.e., intervocalically). This indicates that L1 exemplar patterns, which reflect subphonemic information, are adopted in the L2. Finally, our analysis explains why [z] presents a low

rate of production in L2 English spoken by Brazilian speakers: it is an emerging pattern in the L2, since it has no exemplars from the L1, at least not in word-final position. It will be through experience that these exemplars will become robust and more recurrent.

## Conclusions

The aim of this paper was to investigate [Cs] clusters in English by Brazilian speakers. Its main contribution was to assess the role of orthography, phonological environment and task type not only on the production of epenthetic vowels, but also on the voicing property of the final sibilant. It also considered the role played by the [Cs] ~ [Cis] ongoing sound change from BP into L2 English. Results showed that the orthographic pattern <Ces> favors the production of an epenthetic vowel at higher rates than the <Cs> pattern. As for the task type, it was shown that it was not the visual access to orthographic forms that triggered a vowel to occur, but rather the orthographic patterns.

It was also shown that HNR is strongly influenced by phonetic/phonological environments, including preceding epenthetic vowels, which had not been accounted for in previous studies. We can assume that [z] poses a challenge to Brazilian speakers of L2 English due to the fact that it still has no exemplars from L1 in word-final position.

Concerning the role played by the BP ongoing sound change involving the [Cs] ~ [Cis] alternation, it was shown that robust patterns from the L1 are adopted in L2, including fine phonetic detail that reflects subphonemic properties. This sheds light to the fact that learners not only transfer sounds to the L2, but also phonological behaviors in which such sounds are subject to (as seen with how L2 voicing/HNR is influenced by L1 phonological patterns). We can assume, thus, that better generalizations are posited when the production of [Cs] clusters is assessed globally, rather than accounting for epenthesis and voicing agreement as separate, unrelated phenomena.

## Bibliography

- Paul Boersma. Accurate short-term analysis of the fundamental frequency and the harmonics-to-noise ratio of a sampled sound. In IFA Proceedings 17, pages 97–110, 1993.
- Thais Cristófaró-Silva and Daniela Guimaraes. untitled. Paper presented at Seminário de Ciências da Fala, 9 2021.

1377 Fernanda Delatorre. Brazilian efl learners production of vowel epenthe-  
1378 sis in words ending in-ed. Master's thesis, Florianópolis, SC, 10  
1379 2006.

1380 Carina Silva Fragozo. Aquisição de regras fonológicas do inglês por  
1381 falantes de português brasileiro. PhD thesis, São Paulo, 12 2017.  
1382 URL [http://www.teses.usp.br/teses/disponiveis/8/8139/](http://www.teses.usp.br/teses/disponiveis/8/8139/tde-21122017-124449/)  
1383 [tde-21122017-124449/](http://www.teses.usp.br/teses/disponiveis/8/8139/tde-21122017-124449/).

1384 Silke Hamann and Ilaria E. Colombo. A formal account of the  
1385 interaction of orthography and perception: English intervocalic  
1386 consonants borrowed into Italian. *Natural Language &*  
1387 *Linguistic Theory*, 35(3):683–714, August 2017. ISSN 0167-  
1388 806X, 1573-0859. DOI: 10.1007/s11049-017-9362-3. URL  
1389 <http://link.springer.com/10.1007/s11049-017-9362-3>.

1390 Bruce Hayes. *Introductory Phonology*. 2011. ISBN 9781118315958  
1391 9781444360134. URL [https://nbn-resolving.org/urn:nbn:](https://nbn-resolving.org/urn:nbn:de:101:1-201410212801)  
1392 [de:101:1-201410212801](https://nbn-resolving.org/urn:nbn:de:101:1-201410212801). OCLC: 894709783.

1393 Keith Johnson. *Acoustic and auditory phonetics*. Blackwell Publishers,  
1394 Cambridge, Mass, 1997. ISBN 9780631200949 9780631200956.

1395 Wellington Mendes Jr. and Thais Cristófaros-Silva. Plural formation  
1396 in English: a Brazilian Portuguese case study. In Ubiratã Kickhöfel  
1397 Alves and Jeniffer Imaregna Alcantara de Albuquerque, editors,  
1398 *Second Language Pronunciation: Different Approaches to Teaching*  
1399 *and Training*, volume 64 of *Studies on Language Acquisition*. De  
1400 Gruyter Mouton, 2022 in press.

1401 Yasaman Rafat. The interaction of acoustic and orthographic input  
1402 in the acquisition of Spanish assibilated/fricative rhotics. *Ap-*  
1403 *plied Psycholinguistics*, 36(1):43–66, January 2015. ISSN 0142-  
1404 7164, 1469-1817. DOI: 10.1017/S0142716414000423. URL  
1405 [https://www.cambridge.org/core/product/identifier/](https://www.cambridge.org/core/product/identifier/S0142716414000423/type/journal_article)  
1406 [S0142716414000423/type/journal\\_article](https://www.cambridge.org/core/product/identifier/S0142716414000423/type/journal_article).

1407 Rosane Silveira. O papel desempenhado pelo tipo de tarefa e pela or-  
1408 tografia na produção de consoantes em final de palavra. *Revista de*  
1409 *Estudos da Linguagem*, 15(1), June 2007. DOI: 10.17851/2237-  
1410 2083.15.1.147-180. URL [https://doi.org/10.17851/](https://doi.org/10.17851/2237-2083.15.1.147-180)  
1411 [2237-2083.15.1.147-180](https://doi.org/10.17851/2237-2083.15.1.147-180).

1412 Victor Hugo Medina Soares. Encontros consonantais em final de  
1413 palavra no português brasileiro. Master's thesis, Belo Horizonte,  
1414 2016.



- 1415 Mayara Tsuchida Zanfra. Phonological context as a trigger of voicing  
1416 change: a study on the production of english /s/ and /z/ in word-final  
1417 position by brazilians. Master's thesis, Florianópolis, 12 2013.
- 1418 Chao Zhou. L2 speech learning of European Portuguese /l/ and /flap/  
1419 by L1-Mandarin learners: Experimental evidence and theoretical  
1420 modelling. Language Acquisition, pages 1–2, August 2021. ISSN  
1421 1048-9223, 1532-7817. DOI: 10.1080/10489223.2021.1952591.  
1422 URL [https://www.tandfonline.com/doi/full/10.1080/](https://www.tandfonline.com/doi/full/10.1080/10489223.2021.1952591)  
1423 10489223.2021.1952591.

# Radial Basis Function Artificial Neural Network for Automatic Identification of Interlanguage Transfer Phenomena

ATOS APOLLO SILVA BORGES<sup>1</sup>, BRUNO FERREIRA DE SOUSA<sup>1</sup>,  
ARATUZA RODRIGUES SILVA ROCHA<sup>2</sup>, WILSON JÚNIOR DE ARAÚJO  
CARVALHO<sup>3</sup>, FÁBIO ROCHA BARBOSA<sup>1</sup>, RONALDO MANGUEIRA  
LIMA JÚNIOR<sup>4</sup>

- 1 . Federal University of Piauí
- 2 . Faculdade Afonso Mafrense
- 3 . State University of Ceará
- 4 . Federal University of Ceará

## Abstract

In the recent decades, especially for non-English speaking countries, the modern and more connected world has increased the urgency in learning a second language. Among the obstacles for beginners acquiring a new language are the grapho-phonetic-phonological transfer phenomena between the two language systems, which may undermine their ability to communicate in the target language. The present work proposes a seed for an intelligent software designed to help language learners by providing automatic identification of transfer phenomena produced during their reading process. The algorithm is centered on a Radial Basis Function Artificial Neural Network (RBF-ANN) trained to automatically identify transfer processes between Brazilian Portuguese and English as Foreign Language. Five transfer processes already known in the literature were chosen to demonstrate the concept; however, as an initial approach, the audio samples used for training the algorithm were synthetically generated by the Google Translate™ TTS system. To train the RBF-ANN algorithm we

used the  $f_0$  mean and the mean of the first two Formant Frequencies as signal descriptors. The results presented a promising perspective for the development of a new computer-assisted pronunciation training software (CAPT) with accessible computational resources for Brazilian students and language institutes.

## Introduction

During the learning of a new language, a process called interphonology is manifested. Interphonology is characterized as the creation of a linguistic system different from both the foreign language (L2) and native language (L1), but presenting characteristics from both languages simultaneously [Rocha, 2012]. The students in the processes of acquiring fluency on the second language transfer some of their knowledge of the L1 to the new language due to the already established structure of the L1. This phenomenon, which may be manifested during speech or oral reading, it is called grapho-phonetic-phonological knowledge transfer [Zimmer and Alves, 2006]. The term grapho-phonetic-phonological contemplates not only the transference of phonetic-phonological knowledge but also the transference of the grapheme-phoneme relationship of one language to the other, in the case of this work, Brazilian Portuguese (BP) as L1 to the English as Foreign Language (EFL), the L2. When the learner finds an unknown structure in the foreign language, it uses strategies to adapt L2 to a structure already known in L1.

Transfer phenomena between Portuguese as L1 and English L2 produced by Brazilian learners are well documented in the literature [Silveira et al., 2021]. However, the identification and classification of these processes are made mainly through transcriptions, a slow and laborious process done by specialized linguists. However, there is a shortage of works aimed at recognizing these processes in an automated way. Most studies carry this task as a general mispronunciation identification, comparing the input speech with a pre-recorded dataset of pronunciations, not taking advantage of the nature of each phenomenon. Most studies treat mispronunciation as a random process, not having any pattern or regularity to be explored. Only two works were found proposing forms of automatic identification that take the nature of the phenomena as an important part of the recognition. The first was a categorization of BP speakers by a Self-Organizing Map (SOM) regarding the transfer of stress patterns between BP-L1 and English-L2 [Silva et al., 2011]. The second also aimed to identify transfer processes from BP to English-L2 of Brazilian students using a Multi-Layer Perceptron (MLP) neural network [Rocha, 2017]. The rapid identification of these phenomena would be of great value for

software doing proficiency placement tests and could be used in language schools, distance education, computer-assisted pronunciation training (CAPT), researchers, and inclusion of neurodivergent people [Grund et al., 2020].

Therefore, this work proposes a seed for an intelligent software designed to help language learners by providing automatic identification of transfer phenomena produced during their reading process. The algorithm is centered on a Radial Basis Function Artificial Neural Network (RBF-ANN) trained to automatically identify transfer processes known in the literature of BP transfer to English-L2. The details of the algorithm are described on the RBF Neural Network section. Five transfer processes were chosen to demonstrate the concept and are described on the Acoustic data generation section; however, as an initial approach, the audio samples used for training the algorithm were synthetically generated by the Google Translate™ Text-To-Speech system. We assumed the hypothesis that even simple architectures of Artificial Neural Network, such as Radial Basis Function ANNs, are able to correctly identify the chosen transfer phenomena between Brazilian Portuguese-L1 to English-L2.

## Acoustic data generation

The corpus of this study was constructed using the Corpus of Contemporary American English (COCA), an online and open-access corpus of English with a large variety of written and spoken words. Non-words were also incorporated to the study, all generated by the authors modifying existing words but still obeying English phonological patterns. As the pronunciations in this work should be synthetically generated, there were only two recordings for each word, one with the effects of the transfer phenomenon, as if pronounced by a Brazilian learner, and the other without it, as if pronounced by an English native speaker. A varied quantity of words must be used to be able to reach statistical significance. For this reason, a total of 508 words were used, generating a total of 1016 recordings.

Five widely known transfer phenomena were chosen to be collected in the Google Translate™ TTS system. These phenomena are well documented and commonly found in the pronunciation of Brazilian beginning learners of English.

The first phenomenon investigated was the deletion of initial [h] in words beginning with <h> (henceforth, H-deletion), which corresponds to the deletion of the glottal fricative [h] at the beginning of a word. As initial <h> has no corresponding sound in Portuguese, a

Brazilian learner might produce [i] and [u] in the beginning of ‘hilarious’ and ‘humorist’, respectively.

The second phenomenon was the deletion of initial [h] with a change of [aj] to [i] in words beginning with <hy> (henceforth, HY-i). As in the previous process, the deletion of [h] occurs due to the absence of a sound corresponding to the grapheme <h> in initial position in Portuguese, especially in cognate words such as ‘hyper’, ‘hydrant’ and ‘hydrogen’.

The third process chosen was only changing [aj] to [i] while keeping the pronunciation of initial [h] in words beginning with <hy> (henceforth, HY-hi). The HY-hi process goes in the opposite direction of the previous ones concerning the pronunciation of <h>. In H-deletion and HY-i processes, there is the deletion of initial [h], but in HY-hi the [h] is pronounced, with only a replacement of [aj] by [i], as described above.

The fourth process investigated is the pronunciation of silent <k> with the insertion an epenthetic [i] in words beginning with <kn> (henceforth, KN-kin). This transfer process is characterized by the pronunciation of [k] when <k> should be silent in words like ‘knife’ or ‘knickers’.

The last process investigated was the voicing of /s/ when <s> occurs between two vowels (henceforth, S-z). It is the pronunciation of voiced [z] when it should be voiceless [s]. The voicing occurs in words like ‘basic’, ‘case’ or ‘fantasy’.

After the corpus selection, the speech collection took place on the Google Translate™ online platform. The Text-to-Speech (TTS) system embedded on the platform has the goal of generating a naturally-sounding speech waveform given a text to be synthesized. This process of mapping a sequence of discrete symbols (text) to a real-valued time series (waveform) is design to mimic the human speech production, emulating the periodic and aperiodic components present in human voice.

Recent researchers at Google have proposed the use of neural networks to perform the mapping between linguistic features and acoustic features [Tokuda and Zen, 2016, Zen et al., 2016]. In 2017, about 1/3 of all languages in Google’s TTS options already used Recurrent Neural Networks (RNN) as acoustic models and almost all options of languages in Android mobile devices already used RNN-based TTS systems [Zen, 2017]. The mapping of linguistic features to acoustic features using a parallel-distributed system is remarkably similar to the human reading process in the brain.

To collect the samples produced by Google Translate™, we used the open-source audio software Audacity (version 2.1.2). All the data in this research were collected in August of 2018. The productions were recorded at 44.1 kHz (standard) in Wave 32-bit float PCM. However, raw speech contains thousands of samples, which are often polluted with noise and unnecessary information. The solution is to convert the audio signal into a format with higher information density. To obtain these dense representations, we opted to use the PRAAT software (version 6.0.21).

To test different types of representation, we chose two descriptors: the mean of Formant Frequency (FF) and the mean of the Fundamental Frequency ( $f_0$ ). The PRAAT software presents the oscillogram and spectrogram of audio files. This way, it is possible to select, in each word, the exact region where each researched phenomenon occurred. This specific region was selected, cut, and saved in Wave format, resulting in a file referring to the exclusive region of incidence of transfer processes. The objective was to extract both  $f_0$  mean and the mean of F1 and F2 from the selected region. Although two different methods are used to obtain these values, the same audio file was used for both extractions.

## RBF Neural Network

An artificial neural network is a system composed of ordered neurons in layers interconnected through synaptic weights. These synaptic weights ponder the connection between two neurons, or between an input and a neuron, assuming a higher value according to the influence of that connection to the output of the network. ANN has input nodes that receive stimuli from the external medium and output neurons that provide the network response. Usually, a layer between the input and output neurons is used, known as the hidden layer. The use of the hidden layer structure enables ANN to solve non-linearly separable problems.

A Radial Basis Function Artificial Neural Network is a three-layer feed-forward network that consists of one input layer, responsible for receiving the inputs, one middle layer, fully connected to the input layer, and one output layer, also fully connected by weighted synapses and responsible for outputting the neural network prediction. Each input neuron corresponds to a component of an input vector (in this case F1, F2 and  $f_0$ ). The middle layer consists of N neurons and one bias neuron. Each middle neuron layer neuron computes a kernel

function which is usually the Gaussian function [Hwang and Bang, 1997].

In order to specify the middle layer of an RBF-ANN we have to decide the number of neurons in the kernel layer. The simplest method is creating one neuron for each category present in the data. However, this method is not a good practice for most applications and can be sub-optimal, especially when there is a large number of training patterns. Therefore, we used a K-means algorithm to cluster the training patterns in a reasonable number of groups. K-means is a kind of unsupervised clustering algorithm that search internal groups in the dataset, clustering the samples that belong to the same group enabling the adjustment of the centers and radius of the Gaussian kernels in the hidden layer [Chang et al., 2010].

Next, we use the standard statistical approach to calculate the weights between the middle layer and the output layer. The Least Mean Square Error (LMSE) procedure was used to determine the weights for each synaptic connection between the layers. This method finds the parameters of a linear function by the principle of Least Squares: minimizing the sum of the square of the differences between the observed dependent variable and those predicted by the linear function of the independent variable. In our case, the independent variable is the vector of desired outputs, and the dependent variable is the vector of outputs of the hidden layer. The algorithm finds the linear relationship between these variables (hidden weights) using a nonlinear transformation.

As the neural networks are supervised algorithms, we manually classified the datasets and divided into a training subset (or memory subset) and a testing subset. The training subset is used as reference to the algorithm, presenting enough information about the behavior of the samples to allow for learning and generalization. With the training process completed, the neural network was tested with the testing subset. The samples of the training subset were never presented during the testing process or added in the reference data. This way we could test the accuracy and generalization levels of the model for new samples.

## Results and Discussion

In summary, the results correspond to the average accuracy for each of the 50 iterations using randomized holdout for training, cross-validation and testing subsets. The average accuracy  $\pm$  standard deviation obtained by the algorithm in each phenomenon is distributed in Table 7, presenting the performance in the test sets using both mean  $f_0$  and

the mean of the first two FF. We also displayed the optimal number of hidden neurons found by the K-means algorithm.

Results	Processes			
	H-deletion	HY-i/HY-hi	KN-kin	S-z
Accuracy	0.9203 $\pm$ 0.0234	0.9445 $\pm$ 0.0958	0.9441 $\pm$ 0.0368	0.9308 $\pm$ 0.0223
Kernels	6	2	2	2

The results presented by the RBF-ANN algorithm were in general satisfactory for the identification goal. The algorithm obtained high levels of accuracy for all the phenomena with small variability on the results for the 50 iterations, providing a promising perspective for the development of a new computer-assisted pronunciation training software.

The number of kernels on the hidden layer found by K-means were expected for all process except for H-deletion. One explanation for the higher number of hidden units might be the existences of small clusters on the dataset. These clusters do not only separate the native-like and phenomenon samples, but also reveal internal structures inside the classes. Although these internal clusters are not directly being used to separate the classes globally, they can be used to enhance the accuracy of the decision boundary at the regions of superposition between the classes.

Table 7: Accuracy obtained by the RBF-ANN in each phenomenon studied.

## Conclusions

After the evidence presented by the results, a series of conclusions about the initial hypotheses could be drawn. The results indicated that RBF-ANN can identify the transfer processes produced by the TTS algorithm using the audio descriptor with high levels of accuracy and precision, providing ways to automatically identify the five processes with confidence. The algorithm can be trained with relatively small datasets, and it does not require huge computational power to be trained. These results provided a new perspective on the development of CAPT systems, demonstrating the advantages of using the already developed literature about the transfer phenomena to make the identification process more focused on the transfer patterns. This more efficient approach can be implemented on devices with limited processing power, such as mobile devices and online applications.

For future works, it is still necessary to expand the investigation with more phenomena and to acquire a greater number of samples for each process investigated. Expanding the number samples and testing



new phenomena will provide new information for the development of a simple and efficient identification software. Further investigation can provide significant new information and ideas not only for software development but also about the phenomena themselves.

## Bibliography

Gary W. Chang, Cheng-I Chen, and Yu-Feng Teng. Radial-basis-function-based neural network for harmonic detection. *IEEE Transactions on Industrial Electronics*, 57(6):2171–2179, 2010. DOI: 10.1109/TIE.2009.2034681.

Jonas Grund, Moritz Umfahrer, Lea Buchweitz, James Gay, Arthur Theil, and Oliver Korn. A gamified and adaptive learning system for neurodivergent workers in electronic assembling tasks. In *Proceedings of the Conference on Mensch Und Computer, MuC '20*, page 491–494, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450375405. DOI: 10.1145/3404983.3410420. URL <https://doi.org/10.1145/3404983.3410420>.

Young-Sup Hwang and Sung-Yang Bang. An efficient method to construct a radial basis function neural network classifier. *Neural Networks*, 10(8):1495–1503, 1997. ISSN 0893-6080. DOI: [https://doi.org/10.1016/S0893-6080\(97\)00002-6](https://doi.org/10.1016/S0893-6080(97)00002-6). URL <https://www.sciencedirect.com/science/article/pii/S0893608097000026>.

Aratuza R. S. Rocha. Os efeitos da instrução explícita em fonologia na produção e percepção de consoantes da língua inglesa. Dissertation - Masters, Programa de Pós-Graduação em Linguística Aplicada, Universidade Estadual do Ceará, Fortaleza, Brazil, 2012. URL <http://www.uece.br/posla/dmdocuments/AratuzaRodriguesSilvaRocha.pdf>.

Aratuza Rodrigues Silva Rocha. Identificação de processos de transferência do português do Brasil para o Inglês (L2) por meio de rede neural artificial MLP. PhD, Programa de Pós-Graduação em Linguística Aplicada, Universidade Estadual do Ceará, Fortaleza, Brazil, 2017.

Ana Cristina C. Silva, Ana Cristina P. Macedo, and Guilherme A. Barreto. A SOM-based analysis of early prosodic acquisition of english by brazilian learners: Preliminary results. In *Advances in*

- 1721 Self-Organizing Maps, pages 267–276. Springer Berlin Heidelberg,  
1722 2011.
- 1723 R. Silveira, A. R. Gonçalves, F. Kupske, Ubiratã Kickhöfel Alves, and  
1724 R. M. Lima Jr. Efeito da ortografia. In Investigando os sons de  
1725 línguas não nativas: uma introdução. Editora da Abralín, 2021.
- 1726 K. Tokuda and H. Zen. Directly modeling voiced and unvoiced  
1727 components in speech waveforms by neural networks. In 2016  
1728 IEEE International Conference on Acoustics, Speech and Sig-  
1729 nal Processing (ICASSP), pages 5640–5644, March 2016. DOI:  
1730 10.1109/ICASSP.2016.7472757.
- 1731 Heiga Zen. Generative Model-Based Text-to-Speech Synthesis, 2017.  
1732 URL <https://research.google.com/pubs/pub45882.html>.
- 1733 Heiga Zen, Yannis Agiomyrgiannakis, Niels Egberts, Fergus Hen-  
1734 derson, and Przemysław Szczepaniak. Fast, Compact, and High  
1735 Quality LSTM-RNN Based Statistical Parametric Speech Synthesiz-  
1736 ers for Mobile Devices. San Francisco, CA, USA, June 2016. URL  
1737 <https://research.google.com/pubs/pub45379.html>.
- 1738 Márcia Cristina Zimmer and Ubiratã Kickhöfel Alves. A produção de  
1739 aspectos fonético-fonológicos da segunda língua: instrução explícita  
1740 e conexãoismo. Revista Linguagem & Ensino, 9(2):101–143, 2006.  
1741 ISSN 1983-2400. URL [http://www.rle.ucpel.tche.br/index.  
1742 php/rle/article/view/168](http://www.rle.ucpel.tche.br/index.php/rle/article/view/168).

1743 PART III:

1744 METHODS FOR SPEECH

1745 DATA COLLECTION,

1746 PROCESSING AND ANALYSIS

Fairy tales are more than true: not because  
they tell us that dragons exist, but because  
they tell us dragons can be beaten.

C.K. CHESTERTON

1747

1748

# Behavioral and Neurophysiological Representations of Speech Phonemic Units

ADRIELLE C. SANTANA<sup>1</sup>, ADRIANO V. BARBOSA<sup>2</sup>, HANI C. YEHIA<sup>3</sup>,  
RAFAEL LABOISSIÈRE<sup>4</sup>

1750

1751 1 . Universidade Federal de Ouro Preto

1752 2 . Universidade Federal de Minas Gerais

1753 3 . Université Grenoble Alpes

1754

## Introduction

1755 Many studies showed that we perceive the world around us by cat-  
1756 egorizing the sensory input. This was studied, for example, for the  
1757 perception of emotions [McCullough and Emmorey, 2009] and speech  
1758 [Liberman et al., 1957].

1759 For centuries, researchers around the world try to understand how  
1760 our brain process speech and they try to propose a neurobiological  
1761 model that explains the mechanisms that underlie the production-  
1762 perception relationship. The dual-stream model is one example [Hickok  
1763 and Poeppel, 2007]. But how such models relate or explain the categor-  
1764 ical perception of speech? Many works tried to address this issue.

1765 In the works of Alho et al. [2016], Chevillet et al. [2013] and Möttö-  
1766 nen et al. [2014] the authors worked with a continua based on formant  
1767 variations. In general, in those works the authors concluded about a  
1768 sensorimotor integration of auditory and motor areas for early catego-  
1769 rization which will occur around 120 – 170 *ms* after stimulus onset.  
1770 In Bouton et al. [2018], for a similar continuum, the authors identified  
1771 the encoding of the second formant frequency around 95 – 120 *ms* and  
1772 again at 175 *ms*. In Bidelman et al. [2013] the authors synthesized an

1773 /u/-/a/ continuum and observed categorical perception around 175 *ms*  
 1774 after stimulus onset. With this same continuum Bidelman and Walker  
 1775 [2017] concluded that the phonemic categorization was dependent of  
 1776 attention. However, in Chang et al. [2010] the authors identified the  
 1777 phonemic categorization around 110 *ms* in a task without attention  
 1778 (passive).

1779 Based on those works we propose to perform the investigation of  
 1780 the neural correlates of categorical perception of speech sounds taking  
 1781 into account the attention and the acoustic cue influence as well the  
 1782 brain region measured. We also observed that the works reviewed  
 1783 selected the continuum a priori and did not performed any kind of  
 1784 dissociation of the physical ( $\phi$ ) and psychophysical ( $\psi$ ) characteristics  
 1785 of the stimuli, so we took into account these issues in our study as well.

1786 In order to perform this dissociation we considered the hypothesis  
 1787 illustrated in the Figure 1. It shows four stimuli in a phonemic contin-  
 1788 uum, in this example between the syllables /da/ and /ta/, differing  
 1789 in the Voice Onset Time (VOT). The physical values of VOT is rep-  
 1790 resented in the horizontal axis. The vertical axis shows the probability  
 1791 of /ta/ responses in an identification task. The first (blue) and last  
 1792 (red) stimuli would be unambiguously identified as either /da/ or /ta/.  
 1793 However, the central stimuli (yellow and green), which are close to  
 1794 each other in terms of physical characteristics, would be represented  
 1795 rather distantly from each other in the psychophysical domain. Then,  
 1796 we hypothesize that is possible to identify two separate axes in the neu-  
 1797 rophysiological space: one related to the physical characteristics of the  
 1798 stimuli and another related to its psychophysical categorical perception.

## 1799 Methodology

1800 We performed electroencephalogram (EEG) acquisitions in eleven par-  
 1801 ticipants, right-handed and measured five signals from the electrodes  
 1802 difference: Cz–Tp9, Cz–Tp10, Cz–Fz, Cz–F7 and Cz–F8. The exper-  
 1803 iments were randomized across participants and each one performed  
 1804 both tasks (passive or active) with both continua (based on VOT vari-  
 1805 ations or formants variations). For the acquisitions we selected five  
 1806 stimuli based on the psychometric curve of each subject for each con-  
 1807 tinuum as represented in Figure 2 or a given participant.

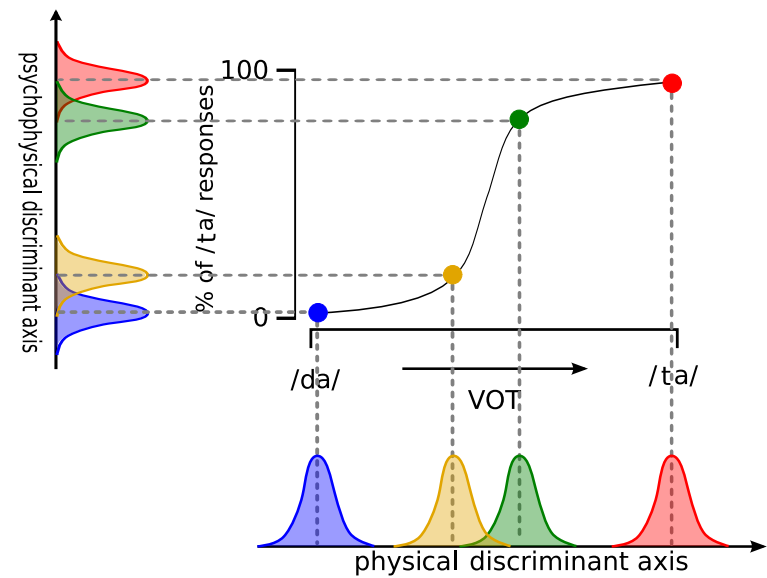


Figure 1: Physical and psychophysical (categorical) neurophysiological axes for the /da/ - /ta/ continuum. The physical values of VOT is represented in the horizontal axis. In the vertical axis, is represented the probability of /ta/ responses in an identification task.

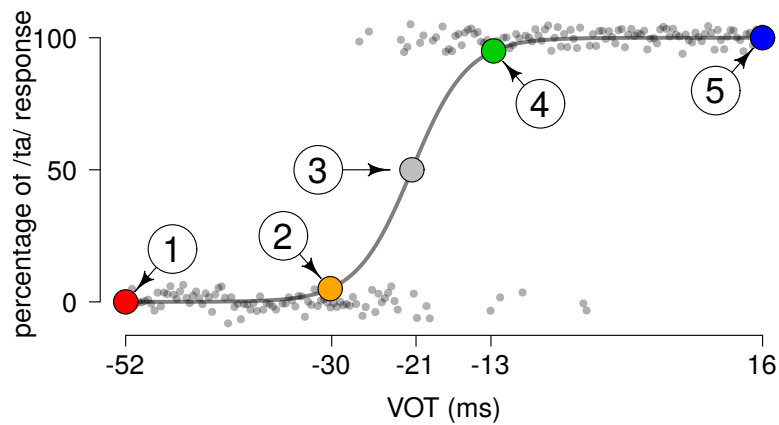


Figure 2: Psychometric curve for a given subject, for the VOT-continuum and the five stimuli selected

## Time-domain processing

To evaluate how the brain oscillations are involved in the coding of the acoustic cues  $e$  to obtain the  $\phi$  and  $\psi$  representations of the stimuli is interesting to work in the time-frequency domain.

For the processing we organized the data by electrode, participant, stimuli, continuum and task. The data was resampled for the execution of the Discrete Wavelet Transform (DWT) with decomposition in 9 levels. This way, the last levels presented a bands similar to those of the main brain oscillations ( $\delta$ ,  $\theta$ ,  $\alpha$ ,  $\beta$  and  $\gamma$ ).

In general, we want to relate the behavior (observed in the psychometric curve) with the neural representation for each participant. However we arrived in a High Dimension Low Sample Size (HDLSS) problem with five observations and 800 wavelet coefficients (features). Then, we developed a regression technique to address this problem named Regression on Low-Dimension Spanned Input Space (RoLD-SIS) [Santana et al., 2020]. Then we were able to obtain the  $\phi$  and  $\psi$  neural discriminant axes which we evaluated in two ways: through the angle between them and through the Euclidean distance between them, which we called discrepancy.

With the angle we compared how it relate with the slope of the psychometric curve, which is reported as a measure of the categorical perception of the participant [Bidelman and Walker, 2017]. With the discrepancy analysis we worked with its value in different regions of the scalogram (graphical representation of the DWT), related to the N1 and P2 latencies and the main brain oscillations as illustrated in Figure 3. We obtained mixed-effects models for each region considering factors related to the continuum, task and electrodes, then we performed an ANOVA followed by a contrast analysis of the factors which presented significant effects.

## Results

### Time-domain results

Following it will be reported our main results for our time-domain analysis and some works that corroborate what we observed.

- We observed that there is a left hemisphere dominance for speech processing [Hickok and Poeppel, 2007, Boemio et al., 2005];
- A spectrotemporal analysis of the acoustic cue happens at the temporal region [Hickok and Poeppel, 2007];

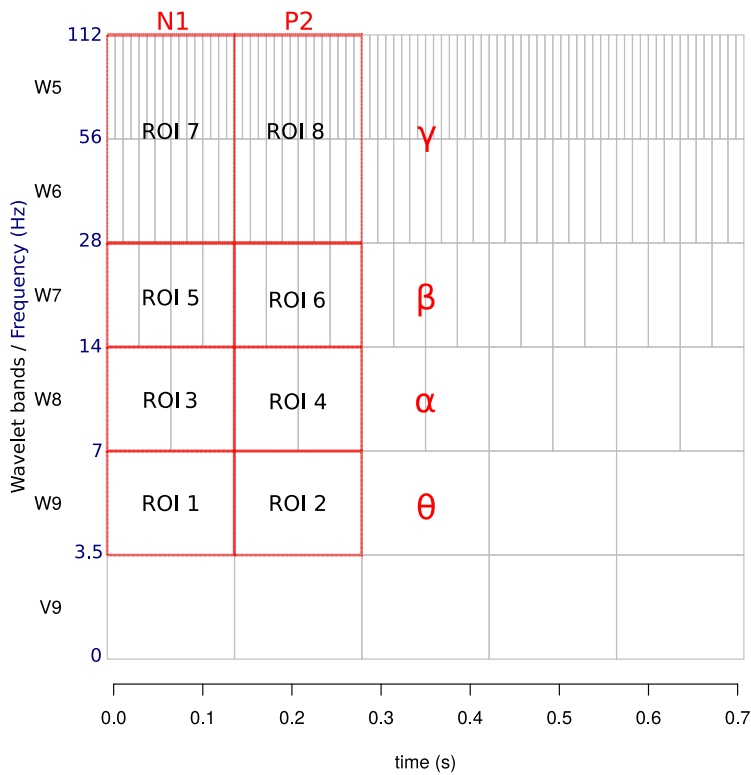


Figure 3: Definition of the ROIs as dependent variables for the models. Each ROI corresponds to specific a frequency band and specific a time interval.



- Generators of N1 and P2 are more laterally localized;
- Formants and VOT evoke different behaviors in N1 and P2 generators;
- N1 is sensitive to VOT variations [Steinschneider et al., 1995, Eggermont, 1995];
- Stimuli are processed differently when there is attention to the task and attention influences the speed of stimuli processing [Möttönen et al., 2014, Alho et al., 2016];
- Ambiguity is reflected into ERP amplitudes [Bidelman and Walker, 2017]
- P2 seems to code ambiguity or effort for speech perception [Rao et al., 2010];
- Attention influences the generators recruited to process stimuli [Hillyard et al., 1973];
- Attention influences more left hemisphere generators than right ones.

#### Time-frequency domain results

Following it will be reported our main results for our time-frequency domain analysis.

- Participants which categorize better have larger difference in internal  $\phi$  and  $\psi$  neural representations of acoustic cues. The Figure 4 illustrates the positive and significant correlation ( $r = 0.788$ ) between the axes' angles and the slopes of the psychometric curve for the formant continuum with the active task;
- It was observed a more lateral location of the speech structures;
- $\gamma$  activity was observed in all scalp regions measured and this is probably related to integration/synchronization of these regions;
- Enhancement of  $\alpha$  activity with attention;
- Larger discrepancies associated with the temporal region than the frontal or medial one;
- Brain oscillations involved in high level speech processing present strong activity as early as the N1 time frame.

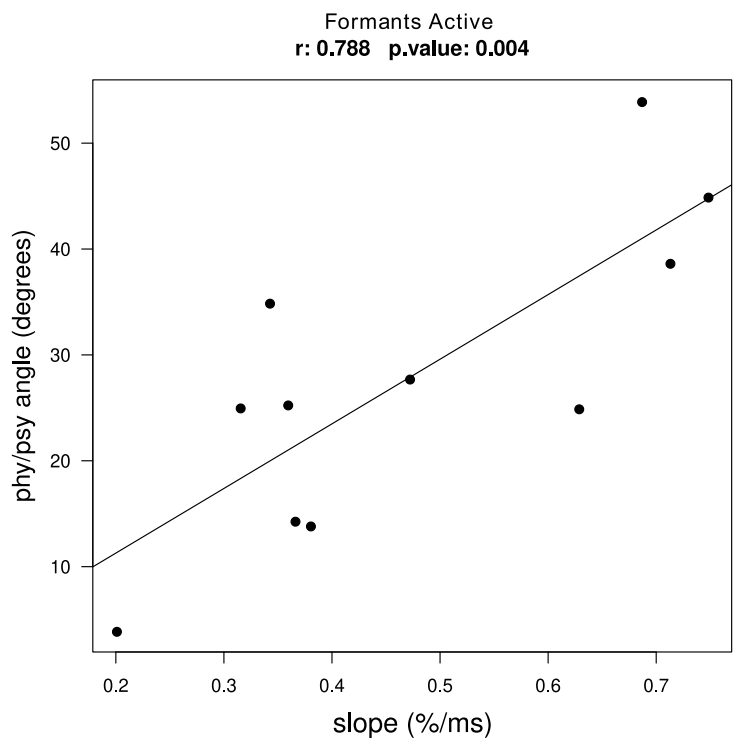


Figure 4: Relationship between the slope of the psychometric curve and the angle between the neurophysiological axes. In this population scatter plot, each point represents a participant. The horizontal and vertical axes represent, respectively, the slope of the fitted psychometric curve at 50% and the angle between the physical and the psychophysical directions obtained by the RoLDSIS procedure. The black line corresponds to the correlation line.

## 1877 Conclusions

1878 In this work we investigated the neural correlates of categorical percep-  
 1879 tion of speech sounds, specifically of Brazilian Portuguese phonemes,  
 1880 evaluating ERPs in the scope of the stimulus acoustic characteristic  
 1881 (VOT and formant frequencies). We studied the brain cortical regions  
 1882 involved in speech perception (temporal and frontal), manipulating the  
 1883 degree of attention to the identification task and using data acquired  
 1884 with the use of a non-invasive method. In our analysis, we propose to  
 1885 identify the physical and psychophysical responses in the ERP, in order  
 1886 to show how the modulations in the time and frequency characteristics  
 1887 of the ERP can be related to the phonemic categorical perception (CP).

1888 We saw that each frequency band and latency seems to code differ-  
 1889 ent aspects of the sound for the speech processing. It was observed that  
 1890 participants who presented behaviorally stronger CP had a larger differ-  
 1891 ence between their physical and psychophysical neural representation  
 1892 of the stimuli. This difference was pronounced for the VOT acoustic  
 1893 cue than for the formants and for active tasks than for the passive ones.  
 1894 It was also shown that the CP occurs when there is no attention to the  
 1895 auditory task but only for the formant-based acoustic cue. Hemispheric  
 1896 differences were observed, with stronger activity at the left hemisphere.  
 1897 Differences were also observed between frontal and temporal cortical  
 1898 regions coded by low-frequency rhythms with more activity at the tem-  
 1899 poral region. In the gamma band we observed no significant difference  
 1900 between the activity at the frontal and temporal regions.

1901 Our results also showed that temporal region structures may also  
 1902 perform some categorization besides the processing of physical acous-  
 1903 tic characteristics of the sounds. We also show how the acoustic cue  
 1904 and task dynamically reconfigure the speech network which should be  
 1905 took into account by a neurobiological model for speech perception.

1906 This study compared different factors related to categorical speech  
 1907 perception in Brazilian Portuguese using a reproducible protocol devel-  
 1908 oped for the study and the evaluation of phonemic categorical percep-  
 1909 tion, and confirmed many of the results found in the literature for other  
 1910 languages.

## 1911 Data and Materials

1912 The data that support the findings of this study, as well as the scripts  
 1913 four reproducing the results and the thesis, are available in the follow-  
 1914 ing repositories:

- 1915 • <https://github.com/Adrielle-Santana/ThesisScripts>

- 1916 • <https://github.com/RoLDSIS/code>
- 1917 • <http://hdl.handle.net/1843/35151>

## 1918 Bibliography

- 1919 Jussi Alho, Brannon M Green, Patrick JC May, Mikko Sams, Hannu  
1920 Tiitinen, Josef P Rauschecker, and Iiro P Jääskeläinen. Early-latency  
1921 categorical speech sound representations in the left inferior frontal  
1922 gyrus. *Neuroimage*, 129:214–223, 2016.
- 1923 Gavin M Bidelman and Brea S Walker. Attentional modulation and  
1924 domain-specificity underlying the neural organization of auditory  
1925 categorical perception. *European Journal of Neuroscience*, 45(5):  
1926 690–699, 2017.
- 1927 Gavin M. Bidelman, Sylvain Moreno, and Claude Alain. Tracing the  
1928 emergence of categorical speech perception in the human auditory  
1929 system. *NeuroImage*, 79:201–212, 2013. ISSN 1053-8119. DOI:  
1930 <https://doi.org/10.1016/j.neuroimage.2013.04.093>.
- 1931 Anthony Boemio, Stephen Fromm, Allen Braun, and David Poeppel.  
1932 Hierarchical and asymmetric temporal sensitivity in human auditory  
1933 cortices. *Nature neuroscience*, 8(3):389, 2005.
- 1934 Sophie Bouton, Valérian Chambon, Rémi Tyrand, Adrian G. Gug-  
1935 gisberg, Margitta Seeck, Sami Karkar, Dimitri van de Ville, and  
1936 Anne-Lise Giraud. Focal versus distributed temporal cortex ac-  
1937 tivity for speech sound category assignment. *Proceedings of*  
1938 *the National Academy of Sciences*, 115(6):E1299–E1308, 2018.  
1939 ISSN 0027-8424. DOI: 10.1073/pnas.1714279115. URL  
1940 <http://www.pnas.org/content/115/6/E1299>.
- 1941 Edward F Chang, Jochem W Rieger, Keith Johnson, Mitchel S Berger,  
1942 Nicholas M Barbaro, and Robert T Knight. Categorical speech repre-  
1943 sentation in human superior temporal gyrus. *Nature neuroscience*, 13  
1944 (11):1428, 2010. DOI: 10.1038/nn.2641.
- 1945 Mark A Chevillet, Xiong Jiang, Josef P Rauschecker, and Maximilian  
1946 Riesenhuber. Automatic phoneme category selectivity in the dorsal  
1947 auditory stream. *Journal of Neuroscience*, 33(12):5208–5215, 2013.
- 1948 Jos J Eggermont. Representation of a voice onset time continuum in  
1949 primary auditory cortex of the cat. *The Journal of the Acoustical*  
1950 *Society of America*, 98(2):911–920, 1995.

- 1951 Gregory Hickok and David Poeppel. The cortical organization of  
 1952 speech processing. *Nature reviews neuroscience*, 8(5):393, 2007.
- 1953 Steven A Hillyard, Robert F Hink, Vincent L Schwent, and Terence W  
 1954 Picton. Electrical signs of selective attention in the human brain.  
 1955 *Science*, 182(4108):177–180, 1973.
- 1956 Alvin M Liberman, Katherine Safford Harris, Howard S Hoffman, and  
 1957 Belver C Griffith. The discrimination of speech sounds within and  
 1958 across phoneme boundaries. *Journal of experimental psychology*, 54  
 1959 (5):358, 1957.
- 1960 Stephen McCullough and Karen Emmorey. Categorical perception  
 1961 of affective and linguistic facial expressions. *Cognition*, 110(2):  
 1962 208–221, 2009.
- 1963 Riikka Möttönen, Gido M van de Ven, and Kate E Watkins. Attention  
 1964 fine-tunes auditory–motor processing of speech sounds. *Journal of*  
 1965 *Neuroscience*, 34(11):4064–4069, 2014.
- 1966 Aparna Rao, Yang Zhang, and Sharon Miller. Selective listening  
 1967 of concurrent auditory stimuli: an event-related potential study.  
 1968 *Hearing research*, 268(1-2):123–132, 2010.
- 1969 Adrielle C. Santana, Adriano V. Barbosa, Yehia Hani C, and Rafael  
 1970 Laboissière. A dimension reduction technique applied to regression  
 1971 on high dimension, low sample size neurophysiological data sets.  
 1972 *BMC Neuroscience*, 2020. (in press).
- 1973 Mitchell Steinschneider, Charles E Schroeder, Joseph C Arezzo, and  
 1974 Herbert G Vaughan. Physiologic correlates of the voice onset time  
 1975 boundary in primary auditory cortex (a1) of the awake monkey:  
 1976 temporal response patterns. *Brain and language*, 48(3):326–340,  
 1977 1995.

1978 PART IV:

1979 BRAZILIAN PORTUGUESE

1980 PHONETICS AND

1981 PHONOLOGY

Fairy tales are more than true: not because  
they tell us that dragons exist, but because  
they tell us dragons can be beaten.

C.K. CHESTERTON

# The implementation of phonic voicing contrast in children's speech: some explorations of clinical data

FABIANA NOGUEIRA GREGIO<sup>1</sup>, ZULEICA CAMARGO<sup>1</sup>  
1 . Pontifícia Universidade Católica de São Paulo

## Abstract

The objective was to investigate the implementation strategies of phonic voicing contrast in Brazilian Portuguese in a group of children with speech disorders in comparison to a control group. From the production of target words with unvoiced and voiced plosive sounds, in contexts of tonic and post-tonic syllables, a set of acoustic measures was extracted and analyzed, a perception experiment was applied, and the acoustic and auditory spheres were explored by means of statistical analysis. The investigation showed that the subjects performed intermediate productions towards the determinant characteristics of the voicing contrast. More than one acoustic cue was implemented for auditory judgment of the voicing contrast.

## Introduction

Speech disorders trigger continual investigations into the refined mechanisms of speech. Among the complaints of speech disorders, the speech therapy clinical setting is faced with the demand to care for cases of "absence/exchange of voiced sounds" [Keske-Soares et al., 2004, Mota et al., 2012]. In traditional phonological views, such disorder is regarded as the absence or substitution of phonemes or dis-

tinctive features. However, clinical evaluations supported by acoustic-phonetic analyses have shown the presence of acoustic cues indicative of the realization of the sound considered absent.

Studies suggest that speakers mark a phonic distinction potentially perceived by them and revealed through intermediate rather than categorical productions, yet not always identified by the listener [Levy, 1993, Ficker, 2003, Gregio and Camargo, 2005, Rodrigues et al., 2008, Gregio et al., 2011, Pereira, 2012].

These productions can be investigated if enlightened by speech production and perception theoretical models that contemplate the dynamic and gradient character of speech, fundamentally relying on the use of speech analysis instruments for explanation [Silva et al., 2001, Gregio et al., 2011, Albano, 2007, Silva, 2010].

Regarding the plosive obstruent consonants of Brazilian Portuguese (BP), the unvoiced-voiced pairs [p]-[b], [t]-[d], and [k]-[g] constitute the repertoire of sounds, characterized by the respective places of articulation bilabial, dental, and velar [Silva et al., 2001, Camargo and Navas, 2008].

Physiologically, the voicing production in plosives involves fine motor coordination of glottic (vocal fold vibration) and supraglottic (vocal tract obstruction) movements [Sweeting and Baken, 1982, Shimizu, 1996, Gregio and Camargo, 2005]. Integrity of lung volume, aerodynamic conditions of the glottis, laryngeal muscles, phono-articulatory organs, and auditory system are required [Hoit et al., 1993, Shimizu, 1996, Hoole et al., 1999].

Acoustically, the production of unvoiced plosives involves the generation of a transient noise source, the acoustic result of complete constriction at some point in the vocal tract followed by its release. In voiced plosives, there are two sound sources: the transient noise coupled with the voice source resulting from the vibration of the vocal folds [Kent and Read, 1992, Johnson, 2003].

One of the acoustic measures used and studied in the investigation of voicing contrasts is voice-onset-time (VOT) [Lisker and Abramson, 1964, Behlau, 1986, Kent and Read, 1992, Levy, 1993, Shimizu, 1996, Cho and Ladefoged, 1999, Camargo et al., 2000, Rocca, 2003]. Other acoustic measures have been reported in voicing contrast studies, such as consonant duration, duration of vowels adjacent to the consonant, fundamental frequency ( $f_0$ ) at the beginning of the vowel following the consonant, frequency of the first formant (F1) at the beginning of the vowel following the consonant, and burst [Barton and Macken, 1980, Shimizu, 1996, Veloso, 1997, van Alphen and Smits, 2004, Benkí,



2046 2005, Lousada et al., 2005, Barroco et al., 2007, Whalen et al., 2007, ?,  
2047 Hanson, 2009, Tachibana et al., 2012].

2048 In the BP literature, especially in the clinical context, VOT is still  
2049 highlighted in the voicing contrast. Studies on this issue in various  
2050 speech situations suggest that more than one acoustic cue seems to be  
2051 involved [Behlau, 1986, Levy, 1993, Barbosa, 1996, Madureira et al.,  
2052 Ficker, 2003, Gregio and Camargo, 2005, Gurgueira, 2006, Barzaghi  
2053 et al., 2007, Bonatto, 2007, de Oliveira e Britto, 2010, Gregio et al.,  
2054 2011, Schliemann, 2011, Souza et al., 2010, Berti et al., 2012, Melo  
2055 et al., 2012, Pereira, 2012].

2056 Thus, the objective of this study was to investigate the implementa-  
2057 tion strategies of phonic voicing contrast in BP in a group of children  
2058 with speech disorders in comparison to a control group.

## 2059 Methods

2060 Six subjects aged between 7 and 10 years old participated in this study,  
2061 two females and four males, of whom three had a diagnosis of speech  
2062 disorders related to voicing contrast (studied group) and three had no  
2063 speech disorders (control group).

2064 The selected subjects presented audiological evaluation with normal  
2065 hearing thresholds, had a negative history of neurological, voice disor-  
2066 ders and other speech disorders unrelated to voicing contrast, and were  
2067 BP speakers with no reference to bilingualism.

2068 The age group did not include the voice change phase, which could  
2069 have affected the voice source as a result of physiological changes  
2070 in the vocal tract [Behlau, 1995]. Within the established age limits,  
2071 differences in laryngeal acoustic measures were not significant between  
2072 male and female children [Andrade, 2009].

2073 The subjects participated in the speech production data collection,  
2074 carried out in a speech laboratory. The corpus consisted of record-  
2075 ings of the subjects reading, with five repetitions in random order,  
2076 sentences following the syllabic structure C1V1C2V2 (consonant1-  
2077 vowel1-consonant2-vowel2). The target words contained the unvoiced  
2078 and voiced plosive sounds of BP (papa, baba, tata, dada, caca, gaga),  
2079 inserted in the carrier-sentence “Diga \_\_\_\_ baixinho”.

2080 The subjects presented different speech productions regarding the  
2081 stress of the target word. Despite the proposed paroxytone stress pat-  
2082 tern (for example, “PApa”), as it is considered the most common in  
2083 BP, some children produced an oxytone stress pattern (“paPA”). Thus,  
2084 the second syllable of the target word was performed as stressed and  
2085 post-tonic. To guarantee the reliability of the data, with the subject

maintaining the same stress of the word throughout its repetitions, the collected corpus was explored through statistical treatment in order to define distinct contexts that could interfere with the reading of the data.

Because the acoustic duration parameter is considered the main correlate of lexical stress in BP (Barbosa, 1996; Aquino, 1997; Gama-Rossi, 1999), duration measures of the V1C2V2 segments were extracted. As a result (figure 1), the acoustic measures showed 100% predictive value in segregating these speech samples into two groups and the most influential variables were duration of C2 and V2, equivalent to the duration of the syllable, finding support in the literature mentioned.

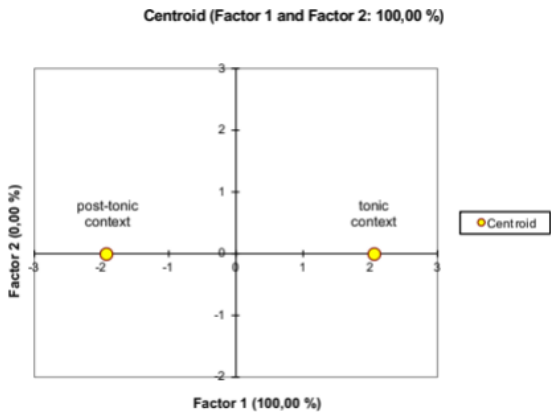


Figure 1: Centroid graph of the discriminant analysis for estimating the subjects' speech productions in the tonic and post-tonic contexts, based on the extracted acoustic measures.

The speech samples were then classified into studied group and control group and according to tonic and post-tonic contexts.

The collected data were analyzed acoustically using the PRAAT software and involved the acoustic inspection of the waveform and broadband spectrogram, and extraction of the measurements: f0 at the beginning of vowel following the consonant (V2); f0 at the stationary point of V2; F1 at the beginning of V2; F1 at the stationary point of V2; measures of duration of the plosive consonant (C2), duration of the previous vowel (V1) and of the following vowel (V2) to the consonant, and duration of the V1C2V2 excerpt of the target word; measures of duration of the VOT; and duration measures of the voicing bar. The voicing bar measures were extracted to contemplate gradient productions and included: duration of the voicing period (voicing bar period in the consonant stretch before the articulation release); duration of the voiceless period (length of the stretch in which there is no voicing bar before the articulation is released); duration of the voicing pre-plosion period (duration of the voicing bar when it is

2114 performed after a period of silence in an excerpt before the articulation  
2115 is released); duration of the total pre-plosion period (total duration of  
2116 the stretch prior to the release of the articulation, regardless of whether  
2117 or not there is a voice bar); plosion duration (burst duration: period  
2118 between the starting point of articulation release and the beginning of  
2119 the vowel). The relative duration measures of all extracted duration  
2120 measures described above were calculated to eliminate influences  
2121 from the subject's speech rate. For sequence of analyses, measures of  
2122 relative duration were used.

2123 The speech production data were submitted to a battery of statis-  
2124 tical tests to consider the existing variants in speech, as proposed in  
2125 research developed by the partnership between researchers and profes-  
2126 sors from LIAAC-PUCSP and the Actuaries and Quantitative Methods  
2127 Department-PUCSP.

2128 The collection of speech perception data was carried out through an  
2129 experiment elaborated using the PRAAT software. The sound stimuli  
2130 consisted of the target words of the carrier sentences of the subjects'  
2131 speech productions. Although the analysis and data reading refer to the  
2132 V1C2V2 excerpt of the target word, for the experiment, the target word  
2133 was edited in its entirety (C1V1C2V2) to avoid the effects of sample  
2134 editing. The stimuli were presented in random order to the judges of  
2135 the perception experiment.

2136 Thirty-nine judges were selected of the same age and level of educa-  
2137 tion, without hearing and/or speech complaints and with no connections  
2138 to the fields of languages and speech-language pathology. The proce-  
2139 dure was performed individually and with the use of headphones, and  
2140 each judge orthographically transcribed the word as he or she heard it  
2141 (word identification).

2142 Next, a statistical analysis was performed to verify the reliability  
2143 of the judges' answers, which resulted in the exclusion of four judges  
2144 (figure 2). To analyze the data from the perception experiment, the  
2145 responses of thirty-five judges were considered.

2146 Based on Johnson [2003], the answers were tabulated in confusion  
2147 matrices. Afterwards, the calculation of the auditory distances of the  
2148 consonant pairs was performed, allowing visualization in the form of  
2149 graphs.

2150 After the speech production data and perception experiment data  
2151 stage, a logistic regression analysis was performed to explore the  
2152 acoustic and auditory spheres. Research was approved by the Ethics  
2153 Committee, protocol number 119/09.

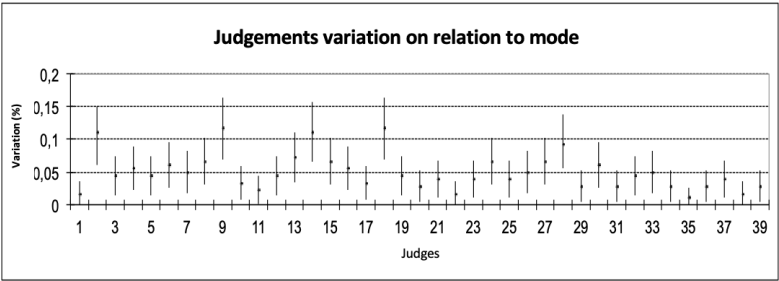


Figure 2: Representation of the auditory judgment variation of each judge in relation to mode value for verification of the judge’s performance.

Results and discussion

Acoustic measures were compared between unvoiced and voiced plosive consonant pairs and, in general, revealed significant differences and allowed for classification of the control and studied groups for both stress contexts (figures 3 and 4).

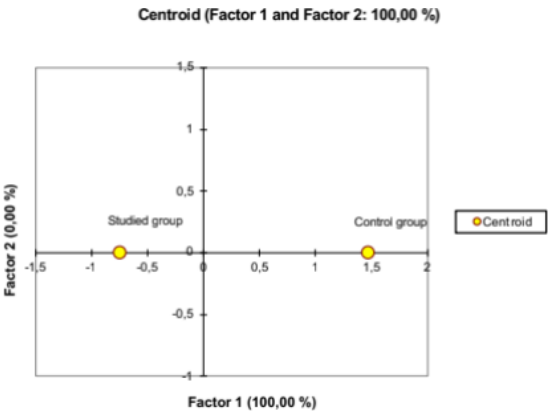


Figure 3: Centroid graph of the discriminant analysis of the estimation of the group of subjects (studied and control) in the tonic context from the acoustic measures.

The results of the calculation of the relative durations of the V1C2V2 excerpt (figures 5 to 8), as well as the voicing bar details (figures 9 to 12), showed differences between the control and studied groups for both stress contexts.

In terms of perception, auditory distances were smaller for samples from the studied group compared to the control group in the tonic context (figure 13). For the post-tonic context, the auditory distances were similar in both groups (figure 14).

The logistic regression analysis revealed that the most influential acoustic measures in the auditory judgments for the unvoiced plosive consonant were, in the tonic context, duration of the previous vowel

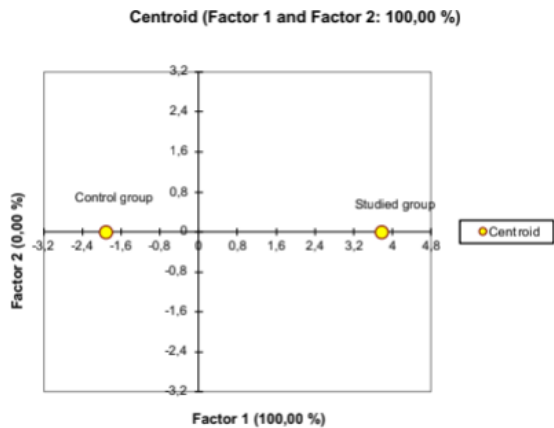


Figure 4: Centroid graph of the discriminant analysis of the estimation of the group of subjects (studied and control) in the post-tonic context from the acoustic measures.

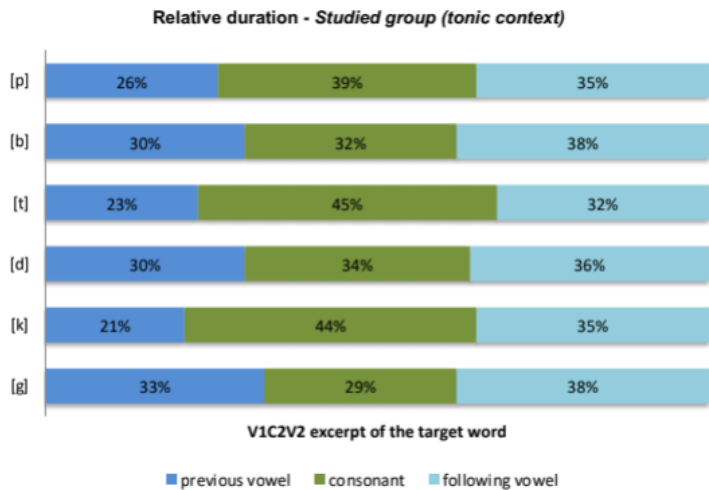


Figure 5: Schematic representation of the relative duration (%) of the vowel preceding the consonant, the plosive consonant and the vowel following the consonant, in relation to the V1C2V2 excerpt of the target word, of the speech samples in the tonic context of the studied group.

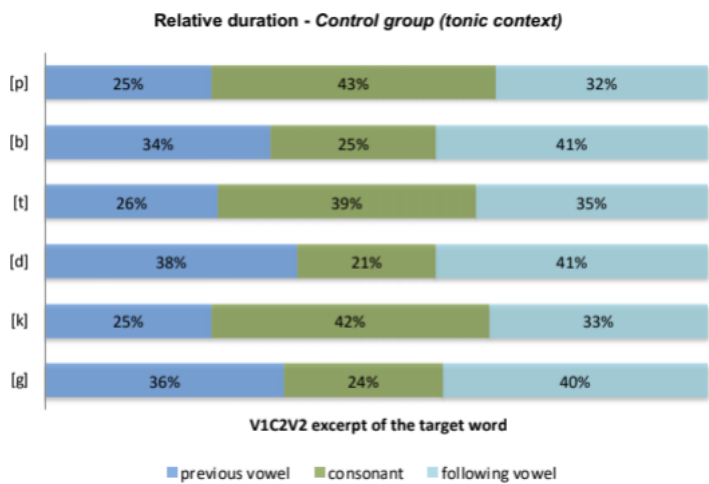


Figure 6: Schematic representation of the relative duration (%) of the vowel preceding the consonant, the plosive consonant, and the vowel following the consonant, in relation to the V1C2V2 excerpt of the target word, of the speech samples in the tonic context of the control group.

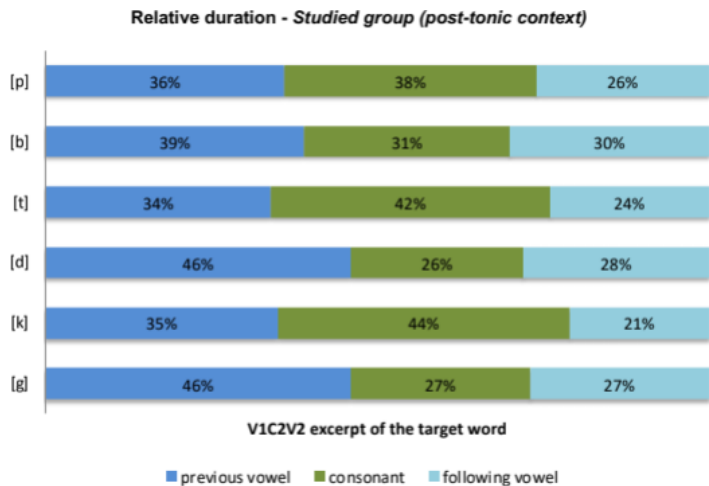


Figure 7: Schematic representation of the relative duration (%) of the vowel preceding the consonant, the plosive consonant, and the vowel following the consonant, in relation to the V1C2V2 excerpt of the target word, of the speech samples in the post-tonic context of the studied group.

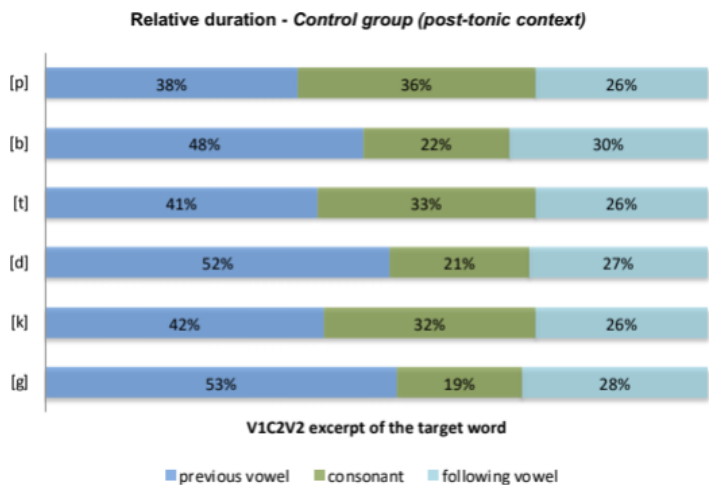


Figure 8: Schematic representation of the relative duration (%) of the vowel preceding the consonant, the plosive consonant, and the vowel following the consonant, in relation to the V1C2V2 excerpt of the target word, of the speech samples in the post-tonic context of the control group.

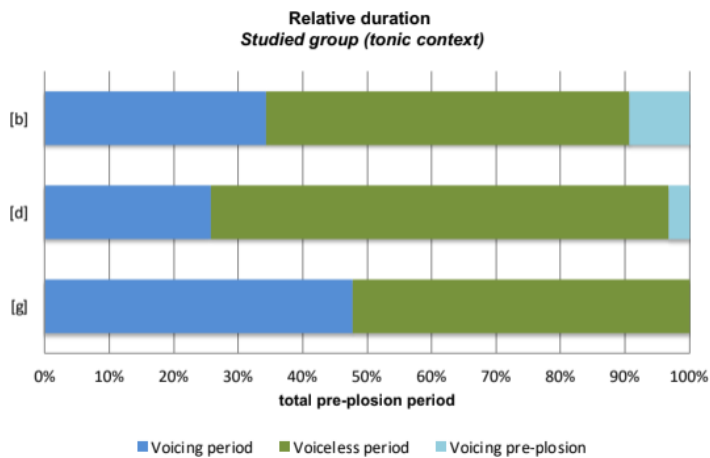


Figure 9: Schematic representation of the relative duration (%) of the voicing period, voiceless period, and the voicing pre-plosion period, in relation to the relative duration of the total pre-plosion period, of the speech samples in the tonic context of the studied group.

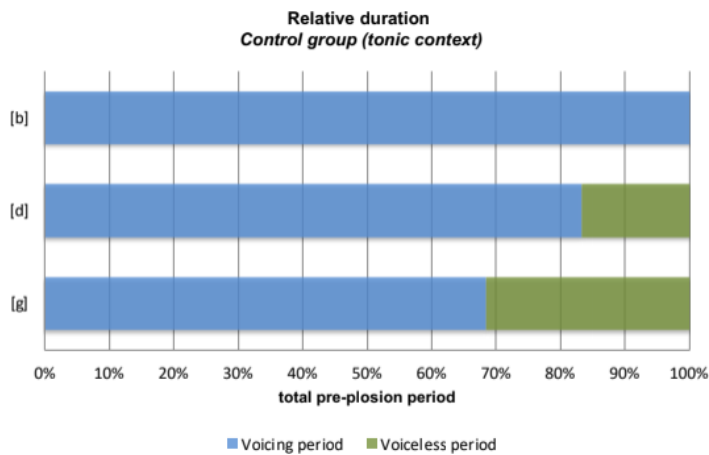


Figure 10: Schematic representation of the relative duration (%) of the voicing period and voiceless period, in relation to the relative duration of the total pre-plosion period, of the speech samples in the tonic context of the control group.

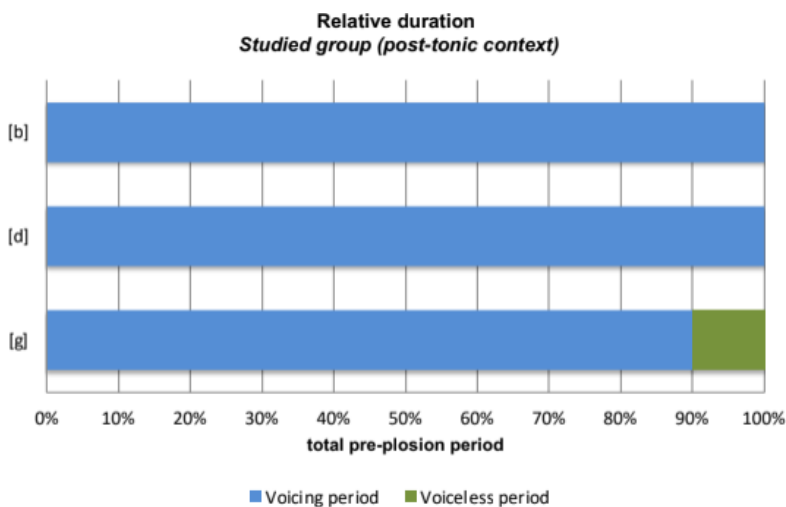


Figure 11: Schematic representation of the relative duration (%) of the voicing period and voiceless period, in relation to the relative duration of the total pre-plosion period, of the speech samples in the post-tonic context of the studied group.



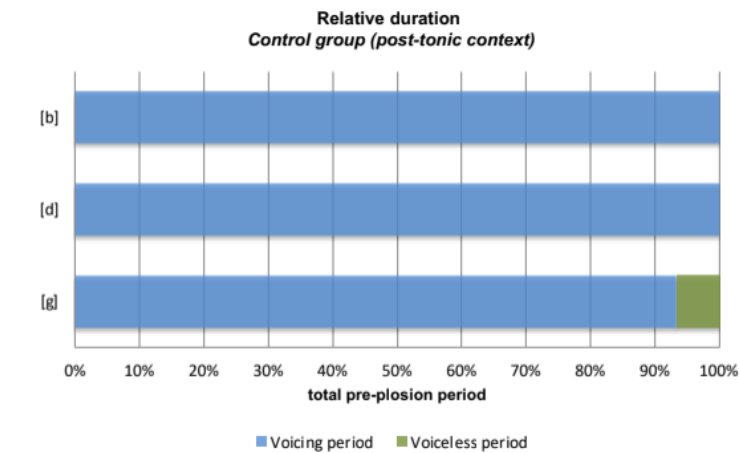


Figure 12: Schematic representation of the relative duration (%) of the voicing period and voiceless period, in relation to the relative duration of the total pre-plosion period, of the speech samples in the post-tonic context of the control group.

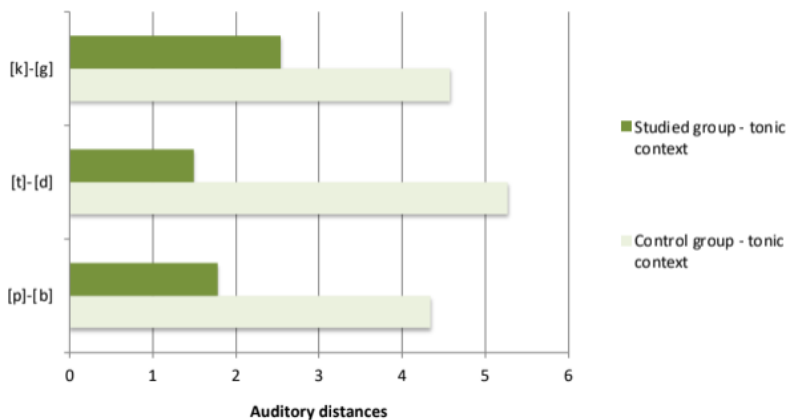


Figure 13: Auditory distances between voiceless and voiced pairs of speech productions in the tonic context of the studied and control groups as a function of the judges’ responses in the perception experiment.

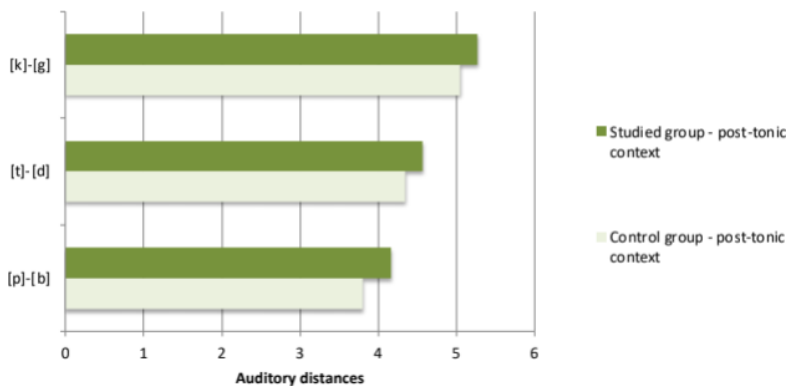


Figure 14: Auditory distances between voiceless and voiced pairs of speech productions in the post-tonic context of studied and control groups as a function of the judges’ responses in the perception experiment.

(V1) and f0 at the beginning of the vowel (V2) (figure 15), and in the post-tonic context, duration of the plosive consonant (C2) and f0 at the beginning of the vowel (V2) (figure 16).

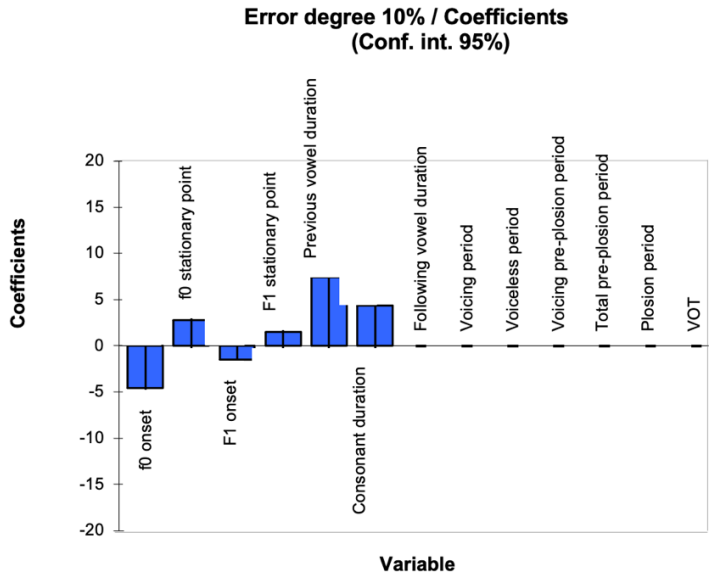


Figure 15: Logistic regression analysis graph for the estimation of auditory judgments from the acoustic measures of the production of unvoiced plosive consonants in the tonic context by the studied and control groups.

For the voiced plosive consonant, the most influential acoustic measures in the auditory judgments were, in the tonic context, the duration of the plosive consonant (C2), the duration of the total pre-plosion and the duration of the voiceless period (figure 17), and in the post-tonic context, the duration of the total pre-plosion period and duration of the voicing period (figure 18).

The VOT duration measure was not revealed as a predictive acoustic cue in the auditory judgment of voicing, in both stress contexts, while other duration measures involved in the voicing bar details were deemed complementary in explaining the implementations that speakers with disorders make and influencing the perception of altered speech.

The data are corroborated by the BP studies that considered the gradient productions from the voicing bar details in the productions of subjects with hearing impairment [Ficker, 2003, Barzaghi et al., 2007, Pereira, 2012] and without speech impairment [Gregio et al., 2011], suggesting the relevance of several acoustic cues for implementing voicing contrast.

In the voiced bilabial plosives of the studied group tonic context, a longer period of voicing pre-plosion period was observed, suggesting

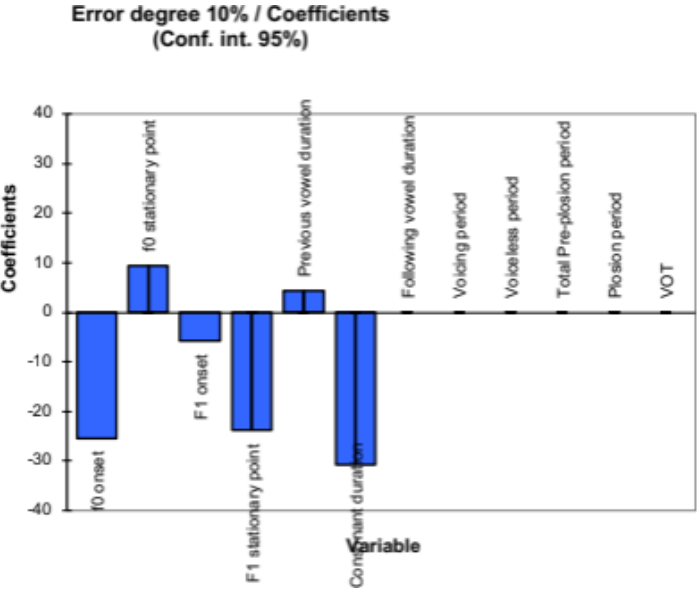


Figure 16: Logistic regression analysis graph for the estimation of auditory judgments from the acoustic measures of the production of unvoiced plosive consonants in the post-tonic context by the studied and control groups.

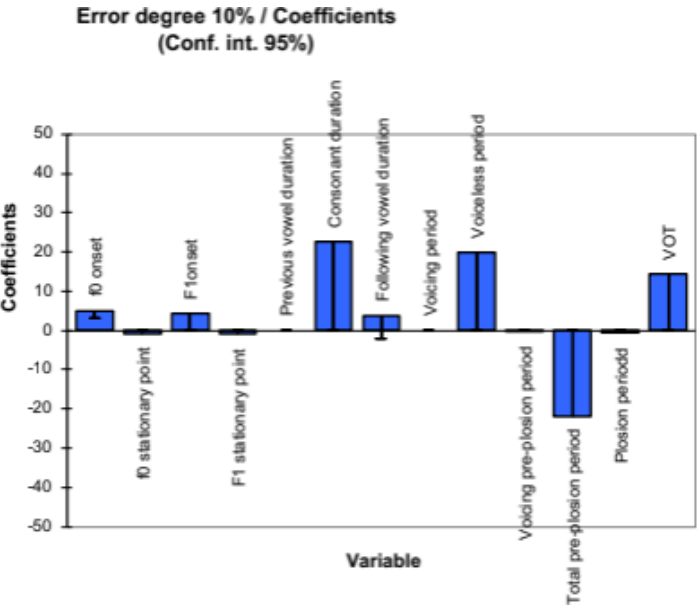


Figure 17: Logistic regression analysis graph for the estimation of auditory judgments from the acoustic measures of the productions of voiced plosive consonants in the tonic context by the studied and control groups.

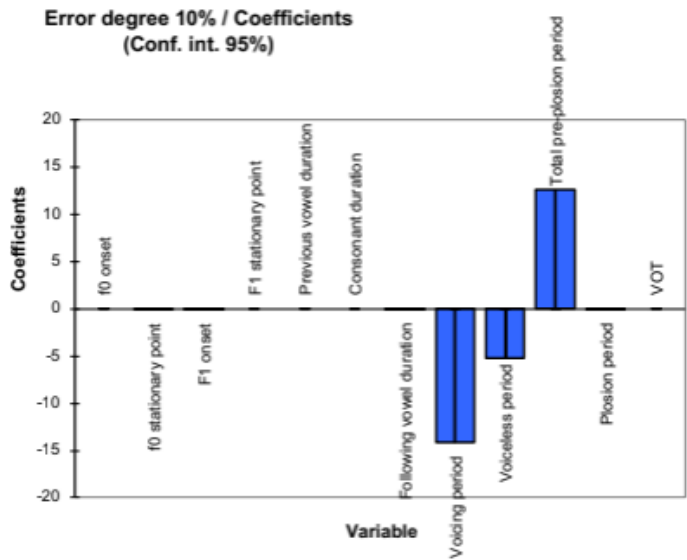


Figure 18: Logistic regression analysis graph for the estimation of auditory judgments from the acoustic measures of voiced plosive consonant productions in the post-tonic context by the studied and control groups.

an attempt to guarantee the necessary voicing period, since bilabials, according to the literature, have a longer voicing period. As for voiced plosive velar, no voicing pre-plosion period was observed, suggesting that the speaker perceives these differentiated cues for the articulatory points in trying to produce the voicing contrast. The literature indicates that velar plosives have a shorter duration of the voicing period [van Alphen and Smits, 2004, Lousada et al., 2005, Barzaghi et al., 2007, Pereira, 2012].

As for the influence of the f0 at the beginning of the vowel, the voicing contrast judgment revealed a predictive value with regard to unvoiced plosive in both stress contexts. As the f0 measurement results from vocal fold activity, which involves aerodynamic and physiological aspects, it tends to be higher at the beginning of the vowel following an unvoiced consonant [Shimizu, 1996]. The f0 measurement has been identified as an acoustic cue in voicing contrasts [Whalen et al., 1993, Hanson, 2009, Gregio et al., 2011]. F1 measures, in turn, did not reveal a predictive value for the auditory judgment of voicing in either stress context.

Regarding the duration measures of the V1C2V2 excerpt, the studied group differentiated the duration of voiced and unvoiced segments, as it kept the voiced plosive consonants' duration shorter than the duration of their respective unvoiced pairs, as expected based on the literature. However, the studied group made this differentiation in a

2216 smaller proportion compared to the control group, suggesting diffi-  
2217 culty in synchronizing the glottic and supraglottic adjustments, given  
2218 different timing of overlapping gestures. The BP literature points to  
2219 higher duration values for vowels preceding and following consonants  
2220 in the production of voiced plosive segments [Barbosa, 1996, Gur-  
2221 gueira, 2006, de Oliveira e Britto, 2010], justifying their influence on  
2222 the auditory judgment of speech disorder.

2223 As a final consideration, the acoustic measures that influenced the  
2224 perception of the auditory judgment of the voicing contrast have been  
2225 shown to be different for each stress context. The listener seems to at-  
2226 tribute different relevance to the acoustic cues involved in the auditory  
2227 judgment of sound. The perception of voicing contrasts showed that  
2228 listeners integrate several acoustic cues to identify and categorize a  
2229 sound. Thus, a clinical diagnosis based on only one acoustic measure-  
2230 ment can be inaccurate.

2231 Most of the speech samples from the studied group, who presented  
2232 clinical demand for speech therapy, could not be categorized as “sound  
2233 exchanges or absences” in view of the data exploration in this study.  
2234 The subjects revealed knowledge about the language, as they per-  
2235 formed intermediate productions towards the determinant charac-  
2236 teristics of the voicing contrast. Such signs denote that the subjects  
2237 perceive differences and seek to implement different actions to support  
2238 the voicing contrast at different articulation points. The acoustic cues  
2239 relevant to the construction of voicing information resided in param-  
2240 eters of duration, which suggest clues about the process of neuromotor  
2241 maturation of speech movements, aspects that have been suggested in  
2242 previous studies with children [Levy, 1993, ?, Albano, 2007].

2243 The demand for temporal refinement surpassed the issues of imple-  
2244 mentation of the f0 acoustic cue. Such aspects relate to the issue of  
2245 synchronization of glottic and supraglottic gestures, which is important  
2246 for the construction of voicing contrast information.

2247 Thus, the exploration of the acoustic signal of the speech samples  
2248 of the studied group indicated an attempt to mark the voicing contrast,  
2249 suggesting that speakers perceive and try to differentiate in their pro-  
2250 duction one sound category from the other, yet these attempts are not  
2251 always processed as relevant information by the listener's perception.

2252 Subjects control and organize their articulatory gestures in terms  
2253 of physical aspects and in terms of perceptual feedback (Gama-Rossi,  
2254 1999; Albano, 2001; Albano, 2007). It is up to the professional to  
2255 guide the child in their attempt to achieve phonic contrast, producing  
2256 articulatory targets that are audible to the listener.

The challenge of working with issues that lie at the interface between speech production and perception offers a rich field of reflection on the nature of speech disorders. Such a challenge may result, in the future, in therapeutic actions that consider the particularities of the manifestation in question, which means contemplating the difficulties and recognizing the implementations made by the speaker, which although not audible at first glance, can be unveiled through instrumental investigation.

## Conclusion

The investigation showed evidence of more than one acoustic cue for the implementation of voicing contrast. The duration of the plosive consonant, voiceless period, and total pre-plosion period (tonic context) and the total pre-plosion period and voicing period (post-tonic context) revealed predictive power of the auditory judgment of the altered speech voicing contrast for voiced plosives. For unvoiced plosive consonants, the influential measures were  $f_0$  at the beginning of the vowel and duration of the previous vowel (tonic context) and  $f_0$  at the beginning of the vowel and duration of the plosive consonant (post-tonic context).

## Bibliography

- Eleonora Cavalcante Albano. Representações dinâmicas e distribuídas: indícios do português brasileiro adulto e infantil. *Letras de Hoje*, 42(1), August 2007. ISSN 1984-7726. URL <https://revistaseletronicas.pucrs.br/index.php/fale/article/view/675>.
- Flávia Viegas de Andrade. Análise de parâmetros espectrais da voz em crianças saudáveis de 4 a 8 anos. Master's thesis, Universidade Veiga de Almeida, Rio de Janeiro, 2009.
- Plínio Almeida Barbosa. At least two macrorhythmic units are necessary for modeling Brazilian Portuguese duration: emphasis on automatic segmental duration generation. *Cadernos de Estudos Linguísticos*, 31, 1996. ISSN 2447-0686. DOI: 10.20396/cel.v31i0.8636991. URL <https://periodicos.sbu.unicamp.br/ojs/index.php/cel/article/view/8636991>.
- Mário André Lopes Barroco, Marta Teresa Pedrosa Domingues, Maria de Fátima Marques de Oliveira Pires, Marisa Lousada, and Luis

- 2293 M. T. Jesus. Análise temporal das oclusivas orais do Português Eu-  
2294 ropeu: um estudo de caso de normalidade e perturbação fonológica.  
2295 Revista CEFAC, 9(2):154–163, June 2007. ISSN 1516-1846. DOI:  
2296 10.1590/S1516-18462007000200003.
- 2297 David Barton and Marlys A. Macken. An instrumental analysis of  
2298 the voicing contrast in word-initial stops in the speech of four-  
2299 year-old english-speaking children. *Language and Speech*, 23  
2300 (2):159–169, April 1980. ISSN 0023-8309, 1756-6053. DOI:  
2301 10.1177/002383098002300203. URL [http://journals.sagepub.](http://journals.sagepub.com/doi/10.1177/002383098002300203)  
2302 [com/doi/10.1177/002383098002300203](http://journals.sagepub.com/doi/10.1177/002383098002300203).
- 2303 Luisa Barzaghi, Kátia Barbosa, and Samar M. El Malt. Deficiência  
2304 de audição e contraste de vozeamento em oclusivas do português  
2305 brasileiro: análise acústica e perceptiva. *Distúrbios da Comunicação*,  
2306 19(3):343–355, 2007.
- 2307 M. S. Behlau. Análise de tempo de início de sonorização na dis-  
2308 criminação dos sons do português. PhD thesis, Escola Paulista de  
2309 Medicina, São Paulo, 1986.
- 2310 Mara Behlau. *Avaliação e tratamento das disfonias*. Editora Lovise,  
2311 São Paulo, 1 edition, January 1995. ISBN 9788585274269.
- 2312 José R. Benkí. Perception of VOT and first formant onset by Spanish  
2313 and English speakers. In *In*, pages 240–248. Cascadilla Press, 2005.
- 2314 Larissa Cristina Berti, Ana Elisa Falavigna, Jéssica Blanca dos Santos,  
2315 and Rita Aparecida de Oliveira. Desempenho perceptivo-auditivo  
2316 de crianças na identificação de contrastes fonológicos entre as  
2317 oclusivas. *Jornal da Sociedade Brasileira de Fonoaudiologia*, 24  
2318 (4):348–354, 2012. ISSN 2179-6491. DOI: 10.1590/S2179-  
2319 64912012000400010.
- 2320 M.T.R.L. Bonatto. *Voices Infantis: a caracterização do contraste do*  
2321 *vozeamento das consoantes plosivas do português brasileiro na fala*  
2322 *de crianças de 3 a 12 anos*. PhD thesis, Pontifícia Universidade  
2323 Católica de São Paulo, São Paulo, 2007.
- 2324 Z. A. Camargo and A. L. G. P. Navas. Fonética e fonologia aplicadas à  
2325 aprendizagem. In J. Zorzi and S. Capellini, editors, *Dislexia e outros*  
2326 *distúrbios da leitura-escrita: letras desafiando a aprendizagem*. Pulso,  
2327 São José dos Campos, 2008.
- 2328 Z.A. Camargo, M. A. S. Fontes, and S. Madureira. *Introdução ao es-*  
2329 *tudo dos sons da fala. apostila da disciplina de Fonética e Fonologia*  
2330 *do curso de Fonoaudiologia-PUCSP*, 2000.

- 2331 Taehong Cho and Peter Ladefoged. Variation and universals in VOT:  
2332 evidence from 18 languages. *Journal of Phonetics*, 27(2):207–229,  
2333 April 1999. ISSN 00954470. DOI: 10.1006/jpho.1999.0094.  
2334 URL [https://linkinghub.elsevier.com/retrieve/pii/](https://linkinghub.elsevier.com/retrieve/pii/S0095447099900943)  
2335 [S0095447099900943](https://linkinghub.elsevier.com/retrieve/pii/S0095447099900943).
- 2336 Ana Teresa Brandão de Oliveira e Britto. Estudo do contraste de  
2337 vozeamento em sujeitos com e sem desvio fonológico. PhD thesis,  
2338 Pontifícia Universidade Católica de Minas Gerais, Belo Horizonte,  
2339 2010.
- 2340 L. Barzaghi Ficker. Estudo da produção e percepção das plosivas do  
2341 português brasileiro por um sujeito com deficiência auditiva. PhD  
2342 thesis, Pontifícia Universidade Católica de São Paulo, São Paulo,  
2343 2003.
- 2344 Fabiana Nogueira Gregio and Zuleica Antonia de Camargo. Dados de  
2345 tempo de início do vozeamento (VOT) na avaliação do sinal vocal  
2346 de indivíduos com paralisia unilateral de prega vocal. *Distúrbios da*  
2347 *Comunicação*, 17(3):289–97, 2005.
- 2348 Fabiana Nogueira Gregio, Renata de Moraes Queiroz, Andrea Baldi  
2349 de Freitas Sacco, and Zuleica Camargo. O uso da eletroglotografia  
2350 na investigação do vozeamento em adultos sem queixa de fala.  
2351 Intercâmbio. *Revista do Programa de Estudos Pós-Graduados em*  
2352 *Linguística Aplicada e Estudos da Linguagem*, 23, 2011. ISSN  
2353 2237-759X. URL [https://revistas.pucsp.br/index.php/](https://revistas.pucsp.br/index.php/intercambio/article/view/8890)  
2354 [intercambio/article/view/8890](https://revistas.pucsp.br/index.php/intercambio/article/view/8890).
- 2355 Adriana Limongeli Gurgueira. Estudo acústico do “voice-onset-time” e  
2356 da duração da vogal na distinção da sonoridade dos sons plosivos em  
2357 crianças com transtorno fonológico. PhD thesis, Universidade de São  
2358 Paulo, São Paulo, 2006.
- 2359 Helen M. Hanson. Effects of obstruent consonants on fundamental  
2360 frequency at vowel onset in English. *The Journal of the Acoustical*  
2361 *Society of America*, 125(1):425–441, January 2009. ISSN 0001-  
2362 4966. DOI: 10.1121/1.3021306. URL [https://www.ncbi.nlm.](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2677272/)  
2363 [nih.gov/pmc/articles/PMC2677272/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2677272/).
- 2364 Jeannette D. Hoit, Nancy Pearl Solomon, and Thomas J. Hixon. Effect  
2365 of lung volume on voice onset time (VOT). *Journal of Speech,*  
2366 *Language, and Hearing Research*, 36(3):516–520, June 1993. ISSN  
2367 1092-4388, 1558-9102. DOI: 10.1044/jshr.3603.516. URL [http:](http://pubs.asha.org/doi/10.1044/jshr.3603.516)  
2368 [//pubs.asha.org/doi/10.1044/jshr.3603.516](http://pubs.asha.org/doi/10.1044/jshr.3603.516).



- 2369 P. Hoole, C. Gobl, and A. N. Chasaide. Laryngeal coarticulation. In  
2370 W. J. Hardcastle and N. Hewlet, editors, *Coarticulation – theory,*  
2371 *data and techniques.* Cambridge University Press, Cambridge, 1999.
- 2372 Keith Johnson. *Acoustic and Auditory Phonetics.* Wiley-Blackwell, 3  
2373 edition, July 2003.
- 2374 Raymond D. Kent and Charles Read. *The acoustic analysis of speech.*  
2375 Singular Pub. Group, San Diego, Calif, 1992. ISBN 9781879105430.
- 2376 Márcia Keske-Soares, Ana Paula Félix Blanco, and Helena Bolli Mota.  
2377 O desvio fonológico caracterizado por índices de substituição e  
2378 omissão. *Revista da Sociedade Brasileira de Fonoaudiologia*, 9(1):  
2379 10–18, 2004.
- 2380 Ivone Panhoca Levy. *Uma nova face da nau dos insensatos: a dificul-*  
2381 *dade de vozear obstruintes em crianças de idade escolar.* PhD thesis,  
2382 Universidade Estadual de Campinas, Campinas, 1993.
- 2383 Leigh Lisker and Arthur S. Abramson. A cross-language study  
2384 of voicing in initial stops: Acoustical measurements. *WORD*,  
2385 20(3):384–422, January 1964. ISSN 0043-7956, 2373-5112.  
2386 DOI: 10.1080/00437956.1964.11659830. URL [http://www.](http://www.tandfonline.com/doi/full/10.1080/00437956.1964.11659830)  
2387 [tandfonline.com/doi/full/10.1080/00437956.1964.](http://www.tandfonline.com/doi/full/10.1080/00437956.1964.11659830)  
2388 [11659830.](http://www.tandfonline.com/doi/full/10.1080/00437956.1964.11659830)
- 2389 M. Lousada, P. Martins, and L. M. T. Jesus. Estudo do pré-  
2390 vozeamento, frequência do burst e locus do f2 das oclusivas orais  
2391 do português europeu. In *Actas do XXI Encontro Nacional da APL*,  
2392 pages 485–494, Porto, Portugal, 2005.
- 2393 Sandra Madureira, Luisa Barzaghi, and Beatriz Mendes. Voicing  
2394 contrasts and the deaf: Production and perception issues. In *Investi-*  
2395 *gations in Clinical Phonetics and Linguistics.*
- 2396 Roberta Michelin Melo, Helena Bolli Mota, Carolina Lisbôa Mez-  
2397 zomo, Brunah de Castro Brasil, Liane Lovatto, and Leonardo  
2398 Arzeno. Parâmetros acústicos do contraste de sonoridade das plosi-  
2399 vas no desenvolvimento fonológico típico e no desviante. *Revista da*  
2400 *Sociedade Brasileira de Fonoaudiologia*, 17(3):304–312, 2012. ISSN  
2401 1516-8034. DOI: 10.1590/S1516-80342012000300012.
- 2402 Helena Bolli Mota, Aline Berticelli, Cintia da Conceição Costa, Fer-  
2403 nanda Marafiga Wiethan, and Roberta Michelin Melo. Ocorrência  
2404 de dessonorização no desvio fonológico: relação com fonemas mais  
2405 acometidos, gravidade do desvio e idade. *Revista da Sociedade*

2406 Brasileira de Fonoaudiologia, 17(4):430–434, December 2012. ISSN  
2407 1516-8034. DOI: 10.1590/S1516-80342012000400011.

2408 Lílian Cristina Kuhn Pereira. As consoantes plosivas do PB: um estudo  
2409 acústico e perceptivo sobre dados de falade sujeitos com deficiência  
2410 auditiva. PhD thesis, Pontifícia Universidade Católica de São Paulo,  
2411 São Paulo, 2012.

2412 Paulina D. Artimonte Rocca. O desempenho de falantes bilíngües:  
2413 evidências advindas da investigação do VOT de oclusivas surdas  
2414 do inglês e do português. DELTA: Documentação de Estudos em  
2415 Linguística Teórica e Aplicada, 19(2):303–328, 2003. ISSN 0102-  
2416 4450. DOI: 10.1590/S0102-44502003000200004.

2417 Luciana Lessa Rodrigues, Maria Cláudia Freitas, Eleonora Cavalcante  
2418 Albano, and Larissa Cristina Berti. Acertos gradientes nos chamados  
2419 erros de pronúncia. Letras, 0(36):85–112, December 2008. ISSN  
2420 2176-1485. DOI: 10.5902/2176148511968. URL [https://](https://periodicos.ufsm.br/letras/article/view/11968)  
2421 [periodicos.ufsm.br/letras/article/view/11968](https://periodicos.ufsm.br/letras/article/view/11968).

2422 Lucila Rey Rocha Schliemann. Contraste de vozeamento por crianças  
2423 entre 6-8 anos: uma abordagem dinâmica. Master’s thesis, Universi-  
2424 dade Estadual de Campinas, Campinas, 2011.

2425 Katsumasa Shimizu. A cross language study of voicing contrasts of  
2426 stop consonants in Asian languages. Seibido, Tokyo, 1996. ISBN  
2427 9784791966486.

2428 Adelaide Silva, Vera Pacheco, and Leonardo Oliveira. Por uma abor-  
2429 dagem dinâmica dos processos fônicos. Revista Letras, 55(0),  
2430 2001. ISSN 2236-0999. DOI: 10.5380/rel.v55i0.2821. URL  
2431 <https://revistas.ufpr.br/letras/article/view/2821>.

2432 Adelaide Hercília Pescatori Silva. O estatuto da análise acústica nos  
2433 estudos fônicos. Cadernos de Letras UFF. Dossiê: Letras e cognição,  
2434 41(1):213–229, 2010.

2435 Ana Paula Ramos de Souza, Lisiane Collares Scott, Carolina Lisbôa  
2436 Mezzomo, Roberta Freitas Dias, and Vanessa Giacchini. Avaliações  
2437 acústica e perceptiva de fala nos processos de dessonorização de  
2438 obstruintes. Revista CEFAC, 13(6):1127–1132, May 2010. ISSN  
2439 1982-0216. DOI: 10.1590/S1516-18462010005000039.

2440 Patricia M. Sweeting and Ronald J. Baken. Voice onset time in a  
2441 normal-aged population. Journal of Speech, Language, and Hearing  
2442 Research, 25(1):129–134, March 1982. ISSN 1092-4388, 1558-9102.

2443 DOI: 10.1044/jshr.2501.129. URL [http://pubs.asha.org/doi/](http://pubs.asha.org/doi/10.1044/jshr.2501.129)  
2444 10.1044/jshr.2501.129.

2445 Ryosuke O. Tachibana, Tatsuya Kitamura, and Masako Fujimoto.  
2446 Differences in articulatory movement between voiced and voiceless  
2447 stop consonants. *Acoustical Science and Technology*, 33(6):391–  
2448 393, 2012. ISSN 1346-3969, 1347-5177. DOI: 10.1250/ast.33.391.  
2449 URL [https://www.jstage.jst.go.jp/article/ast/33/6/33\\_](https://www.jstage.jst.go.jp/article/ast/33/6/33_E1255/_article)  
2450 E1255/\_article.

2451 Petra M. van Alphen and Roel Smits. Acoustical and perceptual  
2452 analysis of the voicing distinction in Dutch initial plosives: the  
2453 role of prevoicing. *Journal of Phonetics*, 32(4):455–491, Octo-  
2454 ber 2004. ISSN 00954470. DOI: 10.1016/j.wocn.2004.05.001.  
2455 URL [https://linkinghub.elsevier.com/retrieve/pii/](https://linkinghub.elsevier.com/retrieve/pii/S0095447004000324)  
2456 S0095447004000324.

2457 João Veloso. Vozeamento, duração e tensão nas oposições de sonori-  
2458 dade das oclusivas orais do português. *Revista da Faculdade de*  
2459 *Letras : Línguas e Literaturas / Linguística*, 14:59–80, 1997.

2460 D. H. Whalen, Arthur S. Abramson, Leigh Lisker, and Maria Mody. F0  
2461 gives voicing information even with unambiguous voice onset times.  
2462 *The Journal of the Acoustical Society of America*, 93(4):2152–2159,  
2463 April 1993. ISSN 0001-4966. DOI: 10.1121/1.406678. URL  
2464 <http://asa.scitation.org/doi/10.1121/1.406678>.

2465 D. H. Whalen, Andrea G. Levitt, and Louis M. Goldstein. VOT in  
2466 the babbling of French- and English-learning infants. *Journal of*  
2467 *phonetics*, 35(3):341–352, July 2007. ISSN 0095-4470. DOI:  
2468 10.1016/j.wocn.2006.10.001. URL [https://www.ncbi.nlm.nih.](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2717044/)  
2469 [gov/pmc/articles/PMC2717044/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2717044/).

2470 PART V:

2471 TEST

Fairy tales are more than true: not because  
they tell us that dragons exist, but because  
they tell us dragons can be beaten.

C.K. CHESTERTON

2472

this is a test