



**HAL**  
open science

# Reducing acoustic overlap of L2 English vowels through gestures encoding lip aperture

Xiaotong Xi, Peng Li, Pilar Prieto

► **To cite this version:**

Xiaotong Xi, Peng Li, Pilar Prieto. Reducing acoustic overlap of L2 English vowels through gestures encoding lip aperture. Alice Henderson; Anastazija Kirkova-Naskova. Proceedings of the 7th International Conference on English Pronunciation: Issues and Practices, pp.261-270, 2023, 10.5281/zenodo.8225191 . hal-04178886

**HAL Id: hal-04178886**

**<https://hal.science/hal-04178886v1>**

Submitted on 8 Aug 2023

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

Xi, X., Li, P., and Prieto, P. (2023). Reducing acoustic overlap of L2 English vowels through gestures encoding lip aperture. In A. Henderson & A. Kirkova-Naskova (Eds.), *Proceedings of the 7<sup>th</sup> International Conference on English Pronunciation: Issues and Practices* (pp. 261–270). Université Grenoble-Alpes. <https://doi.org/10.5281/zenodo.8225191>

## Reducing acoustic overlap of L2 English vowels through gestures encoding lip aperture

Xiaotong Xi  
Pompeu Fabra University, Barcelona

Peng Li  
University of Oslo

Pilar Prieto  
Pompeu Fabra University, Barcelona  
Catalan Institution for Research and Advanced Studies, Barcelona

This study aims to investigate whether hand gestures mimicking the lip aperture of non-native vowels can improve learners' production accuracy after audiovisual perceptual phonetic training. Sixty-six Catalan/Spanish bilingual learners of English were randomly assigned to either the No Gesture or Gesture group for training on the challenging English vowels /æ/ and /ʌ/. In the Gesture group, participants saw the instructor perform gestures mimicking the lip aperture of the low vowel /æ/ and mid vowel /ʌ/. Participants in the No Gesture group only saw the instructor produce the speech. Pronunciation performance was evaluated before and after training using paragraph reading, picture naming, and word imitation tasks. Pillai scores were used to measure the acoustic overlap between /æ/ and /ʌ/. The results showed that although both training groups showed less overlap between the two vowels after training, gestural training had a greater effect than no gesture training in the picture naming and paragraph reading tasks. These findings suggest that hand gestures mimicking visible articulatory features, such as lip aperture, can improve the pronunciation accuracy of non-native sounds in L2 learners.

**Keywords:** hand gestures, English vowels, L2 production, acoustic overlap, phonetic training



---

This chapter is based on the oral presentation given by the authors at the 7th International Conference English Pronunciation: Issues and Practices (EPIP 7) held May 18–20, 2022 at Université Grenoble-Alpes, France. It is licensed under the Creative Commons Attribution 4.0 International License. To view a copy of the license, please go to: <http://creativecommons.org/licenses/by/4.0/>.

## **1 Introduction**

### **1.1 Pronunciation training with multiple modalities**

One of the difficulties adults face as they learn an L2 is the pronunciation of non-native sounds. While these difficulties can be attributed to factors such as the transfer from the L1 phonetic realisation to the L2 production (Kartushina & Frauenfelder, 2014), the limited quantity of the L2 input (Muñoz, 2008), and the quality of the L2 pronunciation instruction (Derwing & Munro, 2015), phonetic training can help overcome these difficulties.

According to the Dual Coding Theory (Paivio, 2014), learners may retain and retrieve information coded through verbal and non-verbal channels more easily than through one modality. Focusing on L2 phonological acquisition, audiovisual phonetic training can boost L2 speech sound learning by providing visual information from facial gestures (e.g., Hazan et al., 2005; Inceoglu, 2016). Furthermore, multimodal phonetic training with embodied cues may provide richer visual sources for L2 speech learning. This is supported by the Embodied Cognition theory, which holds that physical actions can shape cognitive processes (Wilson, 2002). Shapiro and Stolz (2019) as well as Sullivan (2018) argue that embodiment offers benefits in educational settings, where the role of hand gestures in learning has received significant attention. Specifically, while teachers' gestures can help lighten learners' cognitive load by shifting the burden from verbal to visuospatial processing, learners' self-performed gestures indicate whether they have fully comprehended the concept being taught (Shapiro & Stolz, 2019). The benefits of hand gestures have been extensively investigated in different learning domains, such as mathematics, science, and languages, among many others (for a review, see Novack & Goldin-Meadow, 2015). Regarding L2 pronunciation, hand gestures mimicking phonetic features have revealed positive effects on L2 speech production. In the following section, we will review the results of studies investigating the effectiveness of hand gestures in training L2 sound pronunciation.

### **1.2 Phonetic training with hand gestures**

While hand gestures are frequently used by teachers in L2 classrooms to teach pronunciation (Hudson, 2011; Smotrova, 2017), their effectiveness in the learning of L2 sounds has only been tested in a few empirical studies, with mixed results. Some studies have provided positive results of using hand gestures. For example, horizontal sweep gestures mimicking durational features improved the pronunciation accuracy of L2 Japanese long vowels (Li et al., 2020). Similarly, hand gestures mimicking the air burst of stop consonants led to a more accurate pronunciation outcome right after training (Amand & Touhami, 2016; Xi et al., 2020) and at delayed posttest (Li et al., 2021). In contrast, Hoetjes and van Maastricht (2020) used hand gestures to mimic lip rounding and tongue protrusion in order to train Dutch speakers to produce Spanish /u/ and /θ/. However, only the gestures cueing lip rounding of /u/ showed positive effects on pronunciation. Therefore, more evidence is needed to assess the role of hand gestures on the pronunciation of non-native sounds when they mimic the articulatory features of the sounds.

### **1.3 Current study**

The present study examines the effectiveness of multimodal phonetic training with hand gestures encoding articulatory features (specifically, lip aperture) in producing English /æ/ and /ʌ/. Since

Catalan/Spanish speakers lack the /æ/ and /ʌ/ vowels in their native language, they tend to perceive them as a single /a/ category and produce them with smaller spectral differences compared to English speakers (Aliaga-Garcia & Mora, 2008). Formant analysis showed that Catalan/Spanish learners' productions of the /æ/ – /ʌ/ pair were closer to each other compared to native speakers. However, high-variability auditory input training had limited effects on improving their production (Aliaga-Garcia & Mora, 2008). Since better articulatory control is crucial for accurate sound production, including relevant audiovisual and gestural information in phonetic training paradigms could enhance the benefits.

To accurately pronounce the pair of English vowels /æ/ and /ʌ/, two main articulatory features are crucial, namely, tongue height and tongue backness. The first vowel is an open-mid-to-open front vowel, whereas /ʌ/ is an open-mid, central vowel (Carley & Mees, 2020). This indicates that /æ/ is pronounced with a larger lip aperture and a more fronted tongue position than /ʌ/. Following Hudson's (2011) classroom observations, we designed a hand gesture to mirror the lip aperture of the two vowels. The thumb and fingers represent the lower and upper lips, respectively. The vertical distance between them indicates the size of the lip aperture, with a larger distance for /æ/ and a shorter distance for /ʌ/ (see Figure 1).

## Figure 1

*Hand Gestures for the English Vowels /æ/ (Left Panel) and /ʌ/ (Right Panel)*



Based on previous research (e.g., Hoetjes & van Maastricht, 2020; Xi et al., 2020), we hypothesise that training with hand gestures mimicking lip aperture would boost Catalan/Spanish speakers' pronunciation of English /æ/ and /ʌ/ more than training without such gestures.

## 2 Research methodology

### 2.1 Participants

In a between-subjects study with a pretest–posttest–delayed posttest design, we recruited 66 Catalan/Spanish bilinguals (54 females,  $M_{\text{Age}} = 19.7$  years,  $SD = 1.8$ ) from a public university in Catalonia to learn the English vowel pair /æ/ and /ʌ/. They reported having an intermediate English proficiency level and started learning English at an average age of 5.2 years ( $SD = 2.0$ ). None of them reported any hearing or speech impairments. All students volunteered to participate in the experiment and signed a consent form that allowed us to process their data.

Participants were randomly assigned to either the No Gesture (NG) group ( $n = 33$ , female = 26) or the Gesture (G) group ( $n = 33$ , female = 28). Both groups completed the same tests and differed only in the type of training they received. To ensure homogeneity, a series of ANOVA tests confirmed that no significant differences were found in age, age of L2 acquisition, extracurricular hours of English courses, study abroad weeks, and visuospatial working memory between the two groups (all  $p > .05$ ).

## 2.2 Audiovisual materials for the phonetic training

The audiovisual materials for the phonetic training were recorded in a soundproof room. A General American English male speaker acted as the instructor and model speaker. For each training group, we prepared a familiarisation and training video. In the familiarisation video, the instructor provided an explicit explanation in English regarding the differences in lip aperture and tongue position between the two vowels. Especially for the G group, the instructor explained that the hand gesture represented lip aperture, in order to help participants to understand the articulatory feature and map hand gestures to articulatory information. For the training videos, we selected 6 pairs of English CVC minimal word pairs contrasting in /æ/ and /ʌ/ (e.g., *cat-cut*) and created 12 short sentences embedding each of the words (e.g., *A **cat** walks by*). The instructor produced the 12 words and 12 sentences while being video recorded. Again, he performed hand gestures in the G condition when producing the target vowels.

All the video clips were uploaded to the Tobii Pro Lab software<sup>1</sup> to generate the training project for each training group with the familiarisation video preceding the training video. The software allows researchers to easily create multimodal materials (images, videos, webpages, etc.) for carrying out eye-tracking research. The clips of words and sentences were repeated three times. The whole session lasted around 15 minutes for each group.

## 2.3 Procedure

The phonetic training was performed individually in a quiet room. The NG group watched the training video without gestural input during the training, while the G group watched the video with gestural input. Participants did not imitate the speech or hand gestures. They all took the test at three points in time: before, immediately after, and one week after the training with the same tasks. To assess the learning outcome from different angles, we included three tasks: a) a less-controlled paragraph reading task — where participants read an English paragraph<sup>2</sup> containing 14 instances of /æ/ and /ʌ/; b) a spontaneous picture naming task — where participants named 10 simple black-white line drawings<sup>3</sup> designed to elicit 10 words containing /æ/ and /ʌ/; and c) a well-controlled word imitation task — where participants imitated 6 minimal pairs of English CVC words contrasting only in /æ/ and /ʌ/ following the instructor's model speech. While the words in the paragraph reading and picture naming tasks never appeared in the training session, half of the words used in the word imitation task were trained. Participants were audio-recorded during the three testing sessions. After the production tasks, participants did two controlled tasks: a) a language background questionnaire; and b) a symmetry span task to measure the visuospatial

---

<sup>1</sup> <https://www.tobii.com/>

<sup>2</sup> The paragraph was adapted from the textbook *Phonetic Words and Stories*, Book 5 by Kathryn J. Davis. <https://www.soundcityreading.net>

<sup>3</sup> The drawings are from <https://arasaac.org/>

working memory capacity (Blackler et al., 2017). These tasks were included as previous research shows that language experience is associated with L2 pronunciation acquisition (Derwing, 2008), and visuospatial working memory correlates with the learning outcomes through instruction with hand gestures (Aldugom et al., 2020).

### 3 Data analysis and results

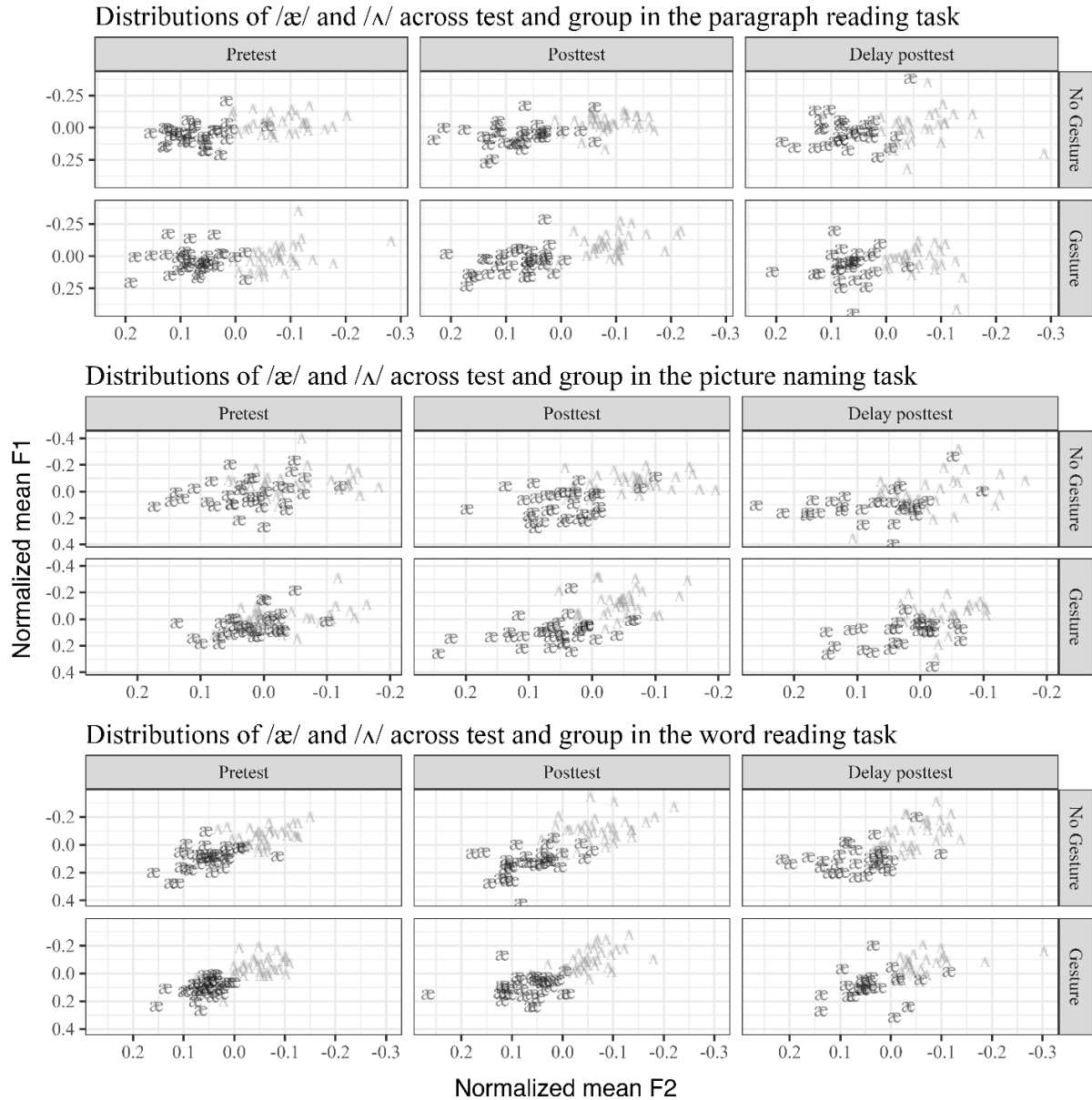
The first author manually annotated 6,794 target vowels using Praat (Boersma & Weenink, 2022). Among these, 129 items containing abnormal formant frequencies and 58 items with mispronunciation (i.e., the participants produced a non-target word) were excluded from the analysis. In addition, eight participants did not do the delayed posttest, and 46 items were not recorded due to technical problems. Thus, 6607 vowels were analysed.

We measured the acoustic distance between the production of the two vowels using Pillai scores. Pillai scores strongly correlate with native listeners' perceptual judgments and are useful in capturing phonetic training gains (Mora, 2021). The Pillai score ranges from 0 to 1, which reflects a speaker's overlap between two vowel realisations (Nycz & Hall-Lew, 2014). A Pillai score closer to 0 means greater acoustic overlap, whereas a score closer to 1 means greater frequency distance. We calculated the Pillai scores of the three tasks for each group at three different testing time points, using the midpoint of the vowels' first and second formants (F1 and F2). The formant values were exported using a script with the maximum number of formants set to 5 and the formant ceiling set to 5500Hz for female speakers and 5000Hz for males. Next, to account for the gender differences in vocal tract length, we normalised the formant values following the Nearey Intrinsic method with phonTools, version 0.2-2.1 (Barreda, 2015), and then averaged the values by the number of testing items. We conducted eighteen Multivariate Analyses of Variance (MANOVA) in RStudio with the normalised mean F1 and F2 mid-point values as the dependent variable and the vowel (2 levels: /æ/ and /ʌ/) as the independent variable. The summary output of each MANOVA gave the Pillai score (see Figure 2).

Table 1 shows the mean Pillai scores across the three tasks at the pretest, posttest, and delayed posttest. For the paragraph reading task, the NG group demonstrated increased overlap immediately and one week after training compared to the pretest. However, the G group showed decreased overlap after training, even though this improvement was not sustained. In the picture naming task, both groups showed decreased vowel overlap after training. The NG group demonstrated an increase in Pillai score at the posttest, followed by a decrease to the pretest level at the delayed posttest. The G group had a larger improvement from the pretest to the posttest, and although their Pillai score decreased after one week, it remained higher than the pretest score. Finally, in the word imitation task, both training groups demonstrated similar improvements in Pillai scores from the pretest to the posttest. However, at the delayed posttest, their Pillai scores decreased to a lower level than the pretest. Overall, the results suggest that hand gestures mimicking articulatory movements have limited effects on improving vowel production in the word imitation task. However, they are effective for improving vowel production in paragraph reading and picture naming tasks.

**Figure 2**

*Vowel F1-F2 Plots with Normalised Mean Formant Values for /æ/ and /ʌ/ across Conditions and Tests from Paragraph Reading (Upper Panel), Picture Naming (Middle Panel), and Word Imitation Tasks (Lower Panel) across Groups and Tests*



**Table 1**

*Pillai Scores across Training Conditions and Testing Time in the Three Production Tasks*

	No Gesture			Gesture		
	Pretest	Posttest	Delayed posttest	Pretest	Posttest	Delayed posttest
Paragraph reading task	0.70	0.69	0.60	0.71	0.74	0.71
Picture naming task	0.32	0.59	0.31	0.17	0.61	0.42
Word imitation task	0.66	0.70	0.53	0.68	0.71	0.52

#### 4 Discussion and conclusion

The present study aimed to investigate the effectiveness of hand gestures that mimic lip aperture on improving the pronunciation of non-native English vowels /æ/ and /ʌ/. A total of 66 Catalan/Spanish learners of English participated in the study, and their pronunciation performance was evaluated through three tasks: a paragraph reading task, a picture naming task, and a word imitation task. The learners' pronunciation improvement was measured using Pillai scores, which are considered a reliable indicator of vowel overlaps in L2 speech production.

The results showed that the effects of hand gestures on vowel production varied depending on the task. In the paragraph reading and picture naming tasks, training with hand gestures reduced the vowel overlap from pretest to posttest more than training without hand gestures. Although after one week the vowel overlap of /æ/ and /ʌ/ became larger (shown by the reduced Pillai score), the overlap of the NG group became larger than the pretest level. In contrast, that of the G group was either back to pretest level in the paragraph reading task or was smaller than the pretest in the picture naming task. These findings suggest that hand gestures can help reduce L2 learners' acoustic overlap of English /æ/ and /ʌ/ in the less-controlled reading task and the spontaneous production task.

However, in the more controlled word imitation task, both training groups showed similar performance, with vowel overlaps of /æ/ and /ʌ/ decreasing immediately after training but increasing beyond the pretest level after one week. This suggests that neither training method helped L2 learners sustain the improvement in the pronunciation of the target vowels. These limited effects of hand gestures on imitating non-native sounds contradict previous studies (e.g., Li et al., 2020; Xi et al., 2020), and the discrepancy could be due to differences in learners' L2 proficiency and task difficulty. While the previous studies trained naïve learners with no experience in the target L2, our participants had intermediate L2 proficiency. Therefore, the controlled imitation task might have been too easy to detect improvements.

Our results are also consistent with a previous phonetic intervention study which identified more benefits from pronunciation intervention in the discourse-reading task than in imitation tasks by English learners with elementary to intermediate levels (Ozakin et al., 2022). In addition, speech production by L2 learners in controlled tasks, such as imitation, may not necessarily reflect their productive knowledge of non-native sounds (Llompert & Reinisch, 2019). Thus, the imitation task may not be an appropriate tool to test the outcome of phonetic training gains in L2 learners with good proficiency levels.



Taken together, the results of the present study suggest that multimodal phonetic training with hand gestures mimicking lip aperture may promote more native-like pronunciation of non-native vowels. Specifically, it helps reduce the acoustic overlap of difficult vowels in spontaneously produced target words and read sentences. Crucially, although participants trained with hand gestures still differed in their production from native English speakers (Pillai score: 0.857, Perry & Tucker, 2019), they showed improvement after a short training session (15 min). Future research could incorporate longer multimodal training sessions to potentially obtain further pronunciation gains. It should be noted that our observation is not based on the significance test, but on the change in Pillai score. As is shown in Table 1, the gestural effects were smaller in the paragraph reading task compared to the picture naming task. This suggests that multimodal training may yield larger effects in spontaneous speech rather than controlled speech production.

To conclude, encoding relevant articulatory information through hand gestures may help learners retrieve information from memory given that the mental representations of the speech sounds are constructed from more than one channel. This interpretation would support the Dual Coding theory (Paivio, 2014). The findings of this study also support the Embodied Cognition Theory (Wilson, 2002) and highlight the importance of involving embodiment in the L2 pronunciation classroom (Sullivan, 2018). Hence, this study supports multimodal pedagogical strategies that encourage teachers to include various sensory modalities in pronunciation instruction (Wrembel, 2011). Hand gestures are effective pedagogical tools that can be easily incorporated into the L2 pronunciation classroom.

## Acknowledgments

This study was supported by the Ministerio de Ciencia, Innovación y Universidades, Agencia Estatal de Investigación and Fondo Europeo de Desarrollo Regional [PGC2018-097007-B-I00], and Agència de Gestió d'Ajuts Universitaris i de Recerca [2017 SGR\_971]. The first author is supported by the Secretaria d'Universitats i Recerca de la Generalitat de Catalunya and the European Social Fund under the Grant for the recruitment of early-stage research staff [2021FI\_B 00137]. The second author is supported by the Research Council of Norway through its Centres of Excellence funding scheme [223265]. The authors sincerely thank Patrick L. Rohrer (Universitat Pompeu Fabra/ Nantes Université) for his help in creating the training materials. Many thanks to the audience at the EPIP7 conference for their useful comments on this study.

## References

- Aldugom, M., Fenn, K., & Cook, S. W. (2020). Gesture during math instruction specifically benefits learners with high visuospatial working memory capacity. *Cognitive Research: Principles and Implications*, 5(27), 1–12. <https://doi.org/10.1186/s41235-020-00215-8>
- Aliaga-Garcia, C., & Mora, J. C. (2008). Assessing the effects of phonetic training on L2 sound perception and production. In M. A. Watkins, A. S. Rauber, & B. O. Baptista (Eds.), *Recent research in second language phonetics/phonology: Perception and production* (pp. 2–31). Cambridge Scholars Publishing.
- Amand, M., & Touhami, Z. (2016). Teaching the pronunciation of sentence final and word boundary stops to French learners of English: Distracted imitation versus audio-visual explanations. *Research in Language*, 14(4), 377–388. <https://doi.org/10.1515/rela-2016-0020>
- Barreda, S. (2015). *phonTools: Functions for phonetics in R* (R package version 0.2-2.1).

- Blacker, K. J., Weisberg, S. M., Newcombe, N. S., & Courtney, S. M. (2017). Keeping track of where we are: Spatial working memory in navigation. *Visual Cognition*, 25(7–8), 691–702. <https://doi.org/10.1080/13506285.2017.1322652>
- Boersma, P., & Weenink, D. (2022). *Praat: doing phonetics by computer* (Version 6.1.51).
- Carley, P., & Mees, I. M. (2020). *American English phonetics and pronunciation practice*. Routledge.
- Derwing, T. M. (2008). Curriculum issues in teaching pronunciation to second language learners. In M. L. Zampini & J. G. Hansen Edwards (Eds.), *Phonology and second language acquisition* (pp. 347–369). John Benjamins Publishing Company.
- Derwing, T. M., & Munro, M. J. (2015). *Pronunciation fundamentals: evidence-based perspectives for L2 teaching and research*. John Benjamins.
- Hazan, V., Sennema, A., Iba, M., & Faulkner, A. (2005). Effect of audiovisual perceptual training on the perception and production of consonants by Japanese learners of English. *Speech Communication*, 47(3), 360–378. <https://doi.org/10.1016/j.specom.2005.04.007>
- Hoetjes, M., & van Maastricht, L. (2020). Using gesture to facilitate L2 phoneme acquisition: the importance of gesture and phoneme complexity. *Frontiers in Psychology*, 11, 575032. <https://doi.org/10.3389/fpsyg.2020.575032>
- Hudson, N. (2011). *Teacher gesture in a post-secondary English as a second language classroom: a sociocultural approach* [Doctoral dissertation, University of Nevada Las Vegas]. UNLV Theses, Dissertations, Professional Papers, and Capstones. <https://doi.org/http://dx.doi.org/10.34917/2432927>
- Inceoglu, S. (2016). Effects of perceptual training on second language vowel perception and production. *Applied Psycholinguistics*, 37(5), 1175–1199. <https://doi.org/10.1017/S0142716415000533>
- Kartushina, N., & Frauenfelder, U. H. (2014). On the effects of L2 perception and of individual differences in L1 production on L2 pronunciation. *Frontiers in Psychology*, 5, 1246. <https://doi.org/10.3389/fpsyg.2014.01246>
- Li, P., Baills, F., & Prieto, P. (2020). Observing and producing durational hand gestures facilitates the pronunciation of novel vowel length contrasts. *Studies in Second Language Acquisition*, 42(5), 1015–1039. <https://doi.org/10.1017/S0272263120000054>
- Li, P., Xi, X., Baills, F., & Prieto, P. (2021). Training non-native aspirated plosives with hand gestures: Learners' gesture performance matters. *Language, Cognition and Neuroscience*, 36(10), 1313–1328. <https://doi.org/10.1080/23273798.2021.1937663>
- Llompарт, M., & Reinisch, E. (2019). Imitation in a second language relies on phonological categories but does not reflect the productive usage of difficult sound contrasts. *Language and Speech*, 62(3), 594–622. <https://doi.org/10.1177/0023830918803978>
- Mora, J. C. (2021). Assessing L2 vowel production gains after high-variability phonetic training: Acoustic measurements vs. perceptual judgements. In J. Romero (Eds.), *Proceedings of the 3rd International Symposium on Applied Phonetics – ISAPh 2021* (pp. 9–18). ISCA. <https://doi.org/10.21437/ISAPh.2021-2>
- Muñoz, C. (2008). Symmetries and asymmetries of age effects in naturalistic and instructed L2 learning. *Applied Linguistics*, 29(4), 578–596. <https://doi.org/10.1093/applin/amm056>
- Novack, M., & Goldin-Meadow, S. (2015). Learning from gesture: How our hands change our minds. *Educational Psychology Review*, 27(3), 405–412. <https://doi.org/10.1007/s10648-015-9325-3>
- Nycz, J., & Hall-Lew, L. (2014). Best practices in measuring vowel merger. *Proceedings of Meetings on Acoustics*, 20(1), 06008. <https://doi.org/10.1121/1.4894063>
- Ozakin, A. S., Xi, X., Li, P., & Prieto, P. (2022). Thanks or Tanks: Training with tactile cues improves learners' accuracy of English interdental consonants in an oral reading task. *Language Learning and Development*, 1–16. Advance online publication. <https://doi.org/10.1080/15475441.2022.2107522>
- Paivio, A. (2014). *Mind and its evolution: A Dual Coding Theoretical Approach*. Psychology Press.
- Perry, S. J., & Tucker, B. V. (2019). L2 production of American English vowels in function words by Spanish L1 speakers. *Canadian Acoustics*, 47(3), 94–95. <https://jcaa.caa-aca.ca/index.php/jcaa/article/view/3328>
- Shapiro, L., & Stolz, S. A. (2019). Embodied cognition and its significance for education. *Theory and*

- Research in Education*, 17(1), 19–39. <https://doi.org/10.1177/1477878518822149>
- Smotrova, T. (2017). Making pronunciation visible: Gesture in teaching pronunciation. *TESOL Quarterly*, 51(1), 59–89. <https://doi.org/10.1002/tesq.276>
- Sullivan, J. V. (2018). Learning and embodied cognition: A review and proposal. *Psychology Learning and Teaching*, 17(2), 128–143. <https://doi.org/10.1177/1475725717752550>
- Wilson, M. (2002). Six views of embodied cognition. *Psychonomic Bulletin & Review*, 9(4), 625–636. <http://view.ncbi.nlm.nih.gov/pubmed/12613670>
- Wrembel, M. (2011). Cross-modal reinforcements in phonetics teaching and learning: An overview of innovative trends in pronunciation pedagogy. In W. S. Lee & E. Zee (Eds.), *Proceedings of the 17th International Congress of Phonetic Sciences* (pp. 104–107). City University of Hong Kong. <http://hdl.handle.net/10593/12093>
- Xi, X., Li, P., Baills, F., & Prieto, P. (2020). Hand gestures facilitate the acquisition of novel phonemic contrasts when they appropriately mimic target phonetic features. *Journal of Speech, Language, and Hearing Research*, 63(11), 3571–3585. [https://doi.org/10.1044/2020\\_JSLHR-20-00084](https://doi.org/10.1044/2020_JSLHR-20-00084)

### About the authors

**Xiaotong Xi** is a PhD student at Universitat Pompeu Fabra in Spain. Her main research interests include second language pronunciation acquisition, phonetic training, and hand gestures.

Emails: [xiaotong.xi@upf.edu](mailto:xiaotong.xi@upf.edu)

**Peng Li** is a postdoc fellow at University of Oslo in Norway. His research focuses on the second language speech learning, multimodal pronunciation training, and prosodic studies.

Emails: [peng.li@iln.uio.no](mailto:peng.li@iln.uio.no)

**Pilar Prieto** is an ICREA research professor at Universitat Pompeu Fabra in Spain. Her main research goal is to understand the role of prosody and co-speech gestures in human communications from a crosslinguistic, developmental, and cognitive perspective.

Emails: [pilar.prieto@upf.edu](mailto:pilar.prieto@upf.edu)