

Modelling English diphthongs with dynamic articulatory targets

Anqi Xu¹, Branislav Gerazov², Daniel van Niekirk¹, Paul Konstantin Krug³,

Santitham Prom-on⁴, Peter Birkholz³, Yi Xu¹

¹Department of Speech Hearing and Phonetic Sciences, University College London, UK

²Faculty of Electrical Engineering and Information Technologies, UCMS, Skopje, RN Macedonia

³Institute of Acoustics and Speech Communication, Technische Universität Dresden, Germany

⁴Computer Engineering Department, King Mongkut's University of Technology Thonburi, Thailand

{a.xu.17, yi.xu}@ucl.ac.uk

Abstract

The nature of English diphthongs has been much disputed. By now, the most influential account argues that diphthongs are phoneme entities rather than vowel combinations. However, mixed results have been reported regarding whether the rate of formant transition is the most reliable attribute in the perception and production of diphthongs. Here, we used computational modelling to explore the underlying forms of diphthongs. We tested the assumption that diphthongs have dynamic articulatory targets by training an articulatory synthesiser with a three-dimensional (3D) vocal tract model to learn English words. An automatic phoneme recogniser was constructed to guide the learning of the diphthongs. Listening experiments by native listeners indicated that the model succeeded in learning highly intelligible diphthongs, providing support for the dynamic target assumption. The modelling approach paves a new way for validating hypotheses of speech perception and production.

Index Terms: diphthongs, computational modelling, articulatory synthesis, American English

1. Introduction

Diphthongs, a special class of vowels, are characterised by transitional formant movements along a path between spectral spaces belonging to two different vowels [1], [2]. Early studies have treated diphthongs as combinations of two vowels, or sometimes vowel-semivowel sequences [3]. However, empirical evidence from a comprehensive work by Gay [4] suggests that English diphthongs are more likely to be distinct phonetic units, based on the observation that listeners were more sensitive to the second formant (F2) movement of the synthetic diphthongs than the onset and offset formant. These results are consistent with more recent findings that the most salient perceptual cue of synthetic diphthongs in noise or reverberation is the intensity of F2 transitions [5]. On the other hand, some studies suggest that the crucial cue in the identification of manipulated diphthongs is the endpoint rather than the transitional trajectory [6]. Another line of studies sought to use a classifier to investigate reliable perceptual cues of diphthongs in a speech corpus, and found that rather than F1-F2 onsets and slopes, classification accuracy was the highest when both F1-F2 onsets and offsets were included [7]. A similar approach was adopted in [8], which reported that incorporating F1-F3 onset, offset and transition rates led to the best classification results.

Not only does the debate about the auditorily relevant formant cues of diphthongs continue, contradictory observations have been made regarding the production of diphthongs. Gay [9] investigated the acoustic properties of five American English diphthongs spoken in three different speech rates from slow to fast. The beginning and terminating vowel formants as well as the rate of F2 movement remained the same across different speaking rates. Further, the final portion of the vowel could be eliminated in fast speech. The unfluctuating formant slopes also accords with more recent acoustic evidence from careful and conversational speech [10], as well as loud speech [11]. As far as articulation is concerned, the tongue body exhibits invariant velocity during the production of diphthongs [12]. Recent X-ray data also show that the tongue flesh points undergo minimal changes in different speaking rates [10]. More importantly, the tongue movements and formant transitions of diphthongs are highly correlated, despite some exceptions [13].

By contrast, some researchers found that spectral changes of diphthongs were lowered in clear speech with prosodic prominence [14]. Unlike previous studies that mainly focus on F2, the measurement of spectral transition is based on the slopes of the first three formants (F1, F2 and F3). They first fitted linear regression lines to the formant slopes and the changes were measured by the root mean-square error of the fitted slopes in different linguistic environments. The inconsistency is probably due to the V-shaped F3 contours of diphthongs [15]. If the gliding movement rather than the onset and offset is the most reliable feature of diphthongs, then one may conclude that diphthongs are distinct phonemes. However, to date there has been little agreement in either perception or production studies.

Here, we used computational modelling to test the assumption that diphthongs have dynamic articulatory targets. We trained an articulatory synthesiser with a 3D vocal tract model to learn real words containing American English diphthongs, following the simulation approach in [16]–[18]. The learning is guided by a phoneme recogniser, comprising a long short-term memory (LSTM) based recurrent neural network to encode a speaker-normalised perceptual space for classifying consonant-to-vowel (CV) sequences. The performance of the articulatory targets are evaluated for the intelligibility of the learned speech in a listening experiment and in terms of the plausibility of the learned articulatory kinematics.

2. Method

2.1. Speech material

Five diphthongs, /aɪ, eɪ, əʊ, aʊ, ɔɪ/, were embedded in real English words with bilabial, alveolar and velar onset consonants, as listed in Table 1. The use of these minimal pairs is to ensure that perception experiments can be carried out naturally by native speakers. Because the two target words ‘bow’ are homographs, we added hints to distinguish the two words, as indicated in brackets. The same hints were also given to the participants during the listening experiment.

Table 1: *Vocal tract parameters in the model.*

Diphthongs	/bV/	/dV/	/gV/
aɪ	buy	die	guy
eɪ	bay	day	gay
əʊ	bow (and arrows)	dough	go
aʊ	(to) bow		
ɔɪ	boy		

2.2. Learning process

We trained a vocal tract model to learn the speech material by an analysis-by-synthesis paradigm as illustrated in Fig. 1. The learning model consists of a production and a perception system. The model begins with exploration of a set of articulatory targets within the parameter range (Fig. 1A). The kinematic trajectories that approach the articulatory targets are based on the timing relations specified by a coarticulation model, which simulates context-sensitive realisation of consonants and vowels (Fig. 1B). The time-varying vocal tract shapes are then converted to cross-sectional area functions for acoustic simulation (Fig. 1C). The synthetic speech is evaluated by a LSTM recurrent neural network that encodes a contrastive phoneme space (Fig. 1D). The model explores the articulatory parameters iteratively, guided by the auditory feedback from the perception system.

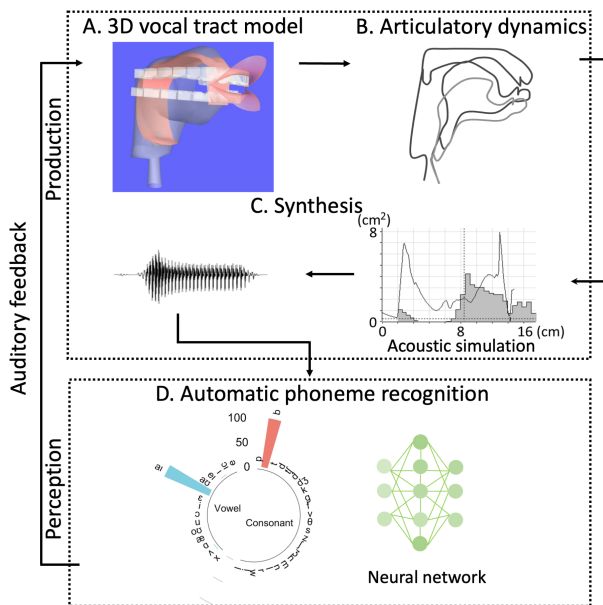


Figure 1: *Overview of the learning process.*

2.3. Vocal tract model

The articulatory synthesiser used in the study is VocalTractLab 2.3 (www.vocaltractlab.de), with a geometrical 3D vocal tract model (Fig. 1A). The vocal tract model was adapted from MRI data of a German male speaker. The synthesiser generates one-dimensional aerodynamic-acoustic simulations based on cross-sectional area functions. The current simulation involved sixteen free vocal tract parameters (Table 2). The vocal folds were set to be fully adducted with moderate tension for the diphthong targets, while the glottis parameters of the consonant targets including the distance between vocal cords, chink area and relative amplitude were free parameters. A falling intonation was added to the synthetic words during the optimisation.

Table 2: *Free vocal tract parameters in the simulation.*

Parameter	Description
HX, HY	Horiz. and vert. hyoid positions
JX, JA	Horiz. jaw position and jaw angle
LP, LD	Lip protrusion and vert. lip distance
TTX, TTY	Horiz. and vert. tongue tip positions
TBX, TBY	Horiz. and vert. tongue blade positions
TCX, TCY	Horiz. and vert. tongue body centre positions
VS	Velum shape
TS1 – TS3	Tongue side elevation from the anterior to the posterior part of the tongue

2.4. Articulatory dynamics

The temporal and spatial movements of the articulators were simulated by a coarticulation model, synchronised dimension-specific sequential target approximation model [16], [19], [20]. In this framework, consonant and vowel articulations are fully synchronised at syllable onset, and despite the consonant-to-vowel (CV) overlap, at the level of individual articulator dimensions, the execution of the articulatory target is sequential. Quantitatively, each articulatory target is represented by height (i.e., positions of the articulators), slope and strength. Unlike monophthongs simulated in [16] with no target slope, the slopes of diphthong targets are free parameters. The coarticulation model generates dynamic trajectories of vocal tract parameters (Fig. 1B) and then the time-series articulatory trajectories will be passed to the articulatory synthesiser for acoustic simulation (Fig. 1C).

2.5. Automatic phoneme recogniser

Previously, we have attempted to use distance metrics with conventional acoustic features such as Mel-frequency cepstral coefficients (MFCCs) [21] to evaluate the synthesised audio and use this to train the vocal tract model but the results were not satisfactory due to difficulties with speaker normalisation [22]. Consequently, we trained a neural network-based phoneme recognition system that learns a speaker-normalised representation (Fig. 1D). We extracted CV sequences with 23 consonants, 11 vowels and 5 diphthongs from the LibriSpeech corpus [23]. To assure that the diphthongs were fully realised, we extracted only CV segments that preceded silences. This resulted with a training set of some 44k segments with a total duration of 4.7h. We applied pre-emphasis (coefficient = 0.97) and calculated the log Mel

spectrogram (25 ms Hamming window, 5ms overlap) with 26 Mel filters (with a maximum frequency of 10 kHz). The log Mel spectrograms were pre-padded to a length of 140 frames (spanning 700 ms). An LSTM recurrent neural network was trained to learn a mapping from the Log Mel spectrograms to a 39-dimensional vector one-hot encoding the CV categories (Fig. 1D). The recogniser had an average of 59% and 95% accuracy in identifying the target consonants and diphthongs respectively for the training set within the CV combinations used in our analysis.

2.6. Optimisation

We use simulated annealing [22] to optimise the vocal tract and glottis parameters. It is a stochastic algorithm that seeks an optimal solution through a coarse-to-fine criterion. This algorithm can heuristically optimise models with many degrees of freedom, such as the speech production system. The learning process started with a neutral position (schwa) followed by adjustments of the vocal tract parameters and gradually converged to a solution. We initiated 20 processes in parallel for each target word, each with 2k iterations. Finally, we manually selected 3 items for each target word to be evaluated in a listening experiment.

2.7. Listening experiment

15 American English native speakers (female: 12; mean age: 35) were invited and screened via Prolific (prolific.co). The learned speech was randomised and presented to the participants via Gorilla, an online experiment tool (gorilla.sc). Before the experiment, the participants filled a brief questionnaire for demographic and language background information. To verify their accents, participants were asked to read the first two sentences of the story “The North Wind and the Sun”, a well-established text recommended by the IPA for eliciting English phonetic contrast. In the experiment, participants were instructed to listen to the audio carefully and choose the word that they heard from the word list. They were allowed to listen to the sounds up to five times. Listeners were asked to undertake the tasks on a computer in a quiet environment without noise or other distractions. Before the identification task, a headphone screening was conducted [23] and five practice trials were presented. The experiment lasted around 12 minutes.

3. Results

Our model learned highly intelligible words containing diphthongs. A demonstration video and learned synthetic samples can be found in https://gitlab.com/Anqi_Xu/dynamic_diphthongs. We calculated the phoneme accuracy based on the response by the native listeners. The phoneme accuracy of the learned consonants and diphthongs across target words is shown in Fig. 2. Error bars show standard errors. The average accuracy was 81.4% and 80.8%, for the consonants and diphthongs, respectively. Bilabial stops had the highest accuracy, followed by alveolar stops and velar stops. With regard to the diphthongs, the ones in ‘day’ and ‘gay’ were perfectly identified and ‘bay’ was also highly intelligible. /aɪ/ had fairly high accuracy in ‘buy’, but the intelligibility was lowered in ‘die’ and ‘guy’. The mean identification accuracy of the diphthongs in all the target

words is summarised in Table 3. The learning of all the diphthongs was fairly successful except /ɔɪ/.

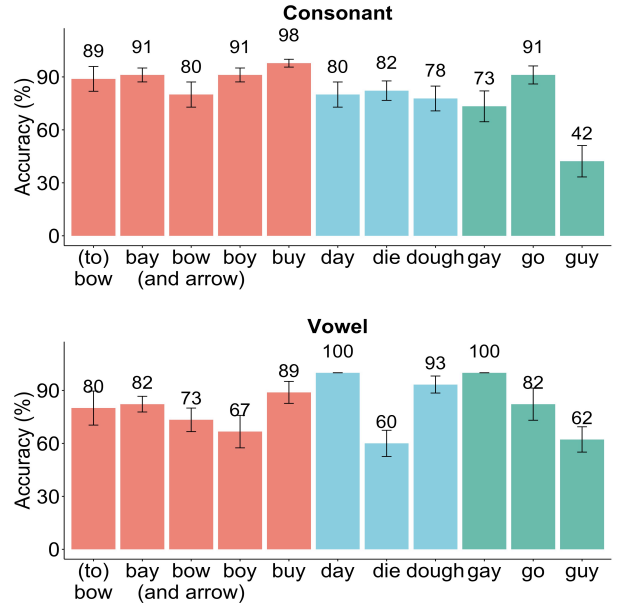


Figure 2: Phoneme accuracy of learned speech in the listening experiment.

Table 3: Mean identification accuracy of diphthongs

/aɪ/	/eɪ/	/əʊ/	/aʊ/	/ɔɪ/
70.4%	94.1%	83.0%	80.0%	66.7%

How listeners classify the synthetic words is shown in Fig. 3. Overall, most of the words were correctly classified. Bilabial stops and alveolar stops were correctly identified most of the time, whereas there was more confusion on the velar stops. For example, ‘gay’ was identified as ‘bay’ and ‘day’ in some cases. With respect to the diphthongs, /aɪ/ and /eɪ/ were often confusing to the listeners. ‘bay’ was sometimes identified as ‘buy’, and ‘die’ as ‘day’. There were a few cases where ‘bow’ (/bəʊ/) was heard as ‘boy’.

Target \ Identified	(to) bow	bay	bow (and arrow)	boy	buy	day	die	dough	gay	go	guy
(to) bow	0	67	7	7	53	13	27	0	0	20	107
bay	7	0	0	0	0	0	7	13	7	233	33
bow (and arrow)	0	53	0	0	0	27	0	0	220	0	0
boy	0	0	7	0	0	7	7	220	0	53	7
buy	0	7	7	0	7	47	167	33	13	13	7
day	0	7	0	0	0	240	0	0	53	0	0
die	13	20	0	0	260	0	0	0	0	0	7
dough	13	13	13	200	33	0	7	0	0	7	13
gay	27	0	180	33	7	7	7	7	0	33	0
go	0	220	0	0	53	13	0	0	13	0	0
guy	240	0	7	13	7	7	0	20	0	7	0

Figure 3: Confusion matrix of the target and identified words in the listening experiment.

Fig. 4 shows the dynamic changes of the spectrograms and the learned vocal tract shapes for bilabial-vowel sequences. The first graph in each row shows the vocal tract shape of the bilabial stops at the moment of maximal constriction. The second and the third graphs show the starting and the ending shapes of the diphthongs. At the syllable onset, although the lips are closed before the release for all the words, the tongue shapes are ready for the dynamic vowel. Take /aɪ/ in synthetic ‘buy’ for example, the initial tongue position is relatively low and later the tongue moves towards a higher position. For /eɪ/, the terminating tongue position is similar to /aɪ/ and /ɔɪ/, while the initial tongue shape seems to be appropriate for a mid vowel. Again, /əʊ/ and /aʊ/ have nearly identical terminating tongue shape, but the initial tongue position is rather different. /aʊ/ starts with a more open vowel shape than /əʊ/. There are also some similarities in the initial tongue positions of /aɪ/ and /aʊ/. Finally, in the case of /ɔɪ/ in ‘boy’, the tongue shape is retracted in the beginning and ends at a higher and more front position.

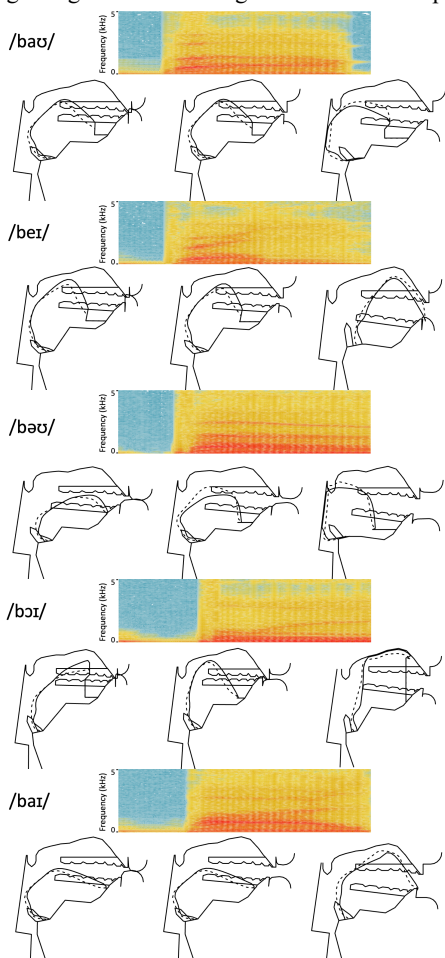


Figure 4: Midsagittal sections of the learned vocal tract shapes and the corresponding spectrograms.

4. Discussion

We have adopted a new approach to probe the nature of diphthongs via computational simulation. We tested the hypothesis that diphthongs are dynamic articulatory targets by training a vocal tract model to learn English diphthongs embedded in real words with the assistance of a phoneme

recogniser. The model learned highly intelligible English words with a mean phoneme accuracy of 81.1% in a multiple choice listening experiment. The results show that the learned dynamic diphthong targets can generate highly intelligible speech. It offers new evidence that diphthongs are likely to be independent phonetic entities with underlying dynamic targets.

The theoretical account of diphthongs as unit phonemes was originally proposed on the basis that formant transition stayed constant in varying speech rates [9], but counterevidence emerged subsequently [14]. We have used a new methodology to address the controversy by emulating the dynamic movement of diphthongs. The learned articulatory targets of diphthongs exhibited beginning and ending vocal tract shapes that resembled two different vowels (Fig. 4). /aʊ/ and /aɪ/ both start with a low and retracted tongue shape; /aɪ/, /eɪ/ and /ɔɪ/ all end with a high and fronted tongue shape; and the ending tongue position of /əʊ/ is similar to that of /aʊ/. The learned vocal tract shapes match well with the tongue positions observed in previous MRI studies [24]. These findings show that diphthongs may have underlying dynamic targets, supporting the proposal of Gay [9].

Another innovation of this study is to use a phoneme recogniser to simulate perceptual guidance, which encodes sound contrasts in a speaker-normalised auditory space. The discrepancies in previous perception research could be attributed to cross-speaker differences in the acoustic realisation of diphthongs. Those studies have identified various auditory signatures of diphthongs, such as F2 transition rates [4], [5], diphthong endpoints [6], F1-F2 onset and offset [7] and all of the above [8]. It is worth noting that some studies use synthetic or manipulated speech [4]–[6], while others are based on the classification results of speech corpus [7], [8]. The classification tasks are dealing with natural speech so the cross-speaker variations may have played a role. The formant onset and offset of diphthongs can be influenced by the anatomical differences between individuals to a great extent [27] and thus they may provide anchoring points for speaker information in the classification tasks. In fact, our pilot study using distance metrics of acoustic features that are not speaker-normalised did not lead to successful learning of diphthongs.

One source of weakness in this study is that the speech data for training the phoneme recogniser is not balanced across all the CV sequences. Thus, the uneven learning performance of the diphthongs could be due to the varied identification accuracy of the recogniser. Moreover, the scope of this study was limited to English diphthongs. Given that there are noticeable cross-linguistic differences in both the perception and production of diphthongs [28], [29], further research should be undertaken to explore how diphthongs in other languages should be modelled. Notwithstanding these limitations, the study directly contributes insights into the dynamic nature of English diphthongs. The computational approach opens a new path towards examining theoretical constructs in speech production and perception.

5. Acknowledgements

This work has been funded by the Leverhulme Trust Research Project Grant RPG-2019-241: "High quality simulation of early vocal learning".

6. References

- [1] I. Lehiste and G. E. Peterson, "Transitions, glides, and diphthongs," *The Journal of the Acoustical Society of America*, vol. 33, no. 3, pp. 268–277, 1961.
- [2] A. Holbrook and G. Fairbanks, "Diphthong formants and their movements," *Journal of Speech and Hearing Research*, vol. 5, no. 1, pp. 38–58, 1962.
- [3] G. L. Trager and H. L. Smith, *An outline of English Structure*. Norman, Oklahoma: Battenburg Press, 1951.
- [4] T. Gay, "A perceptual study of American English diphthongs," *Language and Speech*, vol. 13, no. 2, pp. 65–88, 1970.
- [5] A. K. Nábělek, A. Ovchinnikov, Z. Czyzewski, and H. J. Crowley, "Cues for perception of synthetic and natural diphthongs in either noise or reverberation," *The Journal of the Acoustical Society of America*, vol. 99, no. 3, pp. 1742–1753, 1996.
- [6] A. Bladon, "Diphthongs: A case study of dynamic auditory processing," *Speech Communication*, vol. 4, pp. 145–154, 1985.
- [7] M. Gottfried, J. D. Miller, and D. J. Meyer, "Three approaches to the classification of American English diphthongs," *Journal of Phonetics*, vol. 21, no. 3, pp. 205–229, 1993.
- [8] S. Lee, A. Potamianos, and S. Narayanan, "Developmental acoustic study of American English diphthongs," *The Journal of the Acoustical Society of America*, vol. 136, no. 4, pp. 1880–1894, 2014.
- [9] T. Gay, "Effect of speaking rate on diphthong formant movements," *The Journal of the Acoustical Society of America*, vol. 44, no. 6, pp. 1570–1573, 1968.
- [10] S. M. Tasko and K. Greilick, "Acoustic and articulatory features of diphthong production: A speech clarity study," *Journal of Speech, Language, and Hearing Research*, vol. 53, no. 1, pp. 84–99, 2010.
- [11] K. Tjaden and G. E. Wilding, "Rate and loudness manipulations in dysarthria: acoustic and perceptual findings," *Journal of Speech, Language and Hearing Research*, vol. 47, pp. 766–783, 2004.
- [12] R. D. Kent and K. L. Moll, "Tongue body articulation during vowel and diphthong gestures," *Folia phoniat*, vol. 24, pp. 278–300, 1972.
- [13] C. Dromey, G. O. Jang, and K. Hollis, "Assessing correlations between lingual movements and formants," *Speech Communication*, vol. 55, no. 2, pp. 315–328, 2013.
- [14] J. Wouters and M. W. Macon, "Effects of prosodic factors on spectral dynamics. I. Analysis," *The Journal of the Acoustical Society of America*, vol. 111, no. 1, pp. 417–427, 2002.
- [15] F. Clermont, "Spectro-temporal description of diphthongs in F1-F2-F3 space," *Speech Communication*, vol. 13, pp. 377–390, 1993.
- [16] A. Xu, P. Birkholz, and Y. Xu, "Coarticulation as synchronized dimension-specific sequential target approximation: An articulatory synthesis simulation," *International Congress of Phonetic Sciences ICPhS*, 2019.
- [17] D. R. van Niekirk, A. Xu, B. Gerazov, P. K. Krug, P. Birkholz, and Y. Xu, "Finding Intelligible Consonant-vowel sounds using high-quality articulatory synthesis," in *Interspeech*, 2020, pp. 4457–4461.
- [18] S. Prom-On, P. Birkholz, and Y. Xu, "Identifying underlying articulatory targets of Thai vowels from acoustic data based on an analysis-by-synthesis approach," *EURASIP Journal on Audio, Speech, and Music Processing*, 23, 2014.
- [19] Y. Xu, "Syllable as a synchronization mechanism that makes human speech possible," *PsyArXiv*. 2020.
- [20] Zirui. Liu, Yi. Xu, and F. Hsieh, "Coarticulation as synchronised CV co-onset – Parallel evidence from articulation and acoustics," *Journal of Phonetics*, vol. 90, p. 101116, 2022.
- [21] S. B. Davis and P. Mermelstein, "Comparison of parametric representations for monosyllabic word recognition in continuously spoken sentences," *IEEE Transactions on Acoustics, Speech and Signal Processing*, vol. 28, no. 4, pp. 357–366, 1980.
- [22] K. Johnson and M. J. Sjerps, "Speaker normalization in speech perception," in *The handbook of speech perception*, J. S. Pardo, L. C. Nygaard, R. E. Remez, and D. B. Pisoni, Eds. Wiley, 2021.
- [23] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: An ASR corpus based on public domain audio books," in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2015, pp. 5206–5210.
- [24] S. Kirkpatrick, C. D. Gelatt, and M. P. Vecchi, "Optimization by simulated annealing," *Science*, vol. 220, no. 4598, pp. 671–680, 1983.
- [25] K. J. P. Woods, M. H. Siegel, J. Traer, and J. H. McDermott, "Headphone screening to facilitate web-based auditory experiments," *Attention, Perception, and Psychophysics*, vol. 79, no. 7, pp. 2064–2072, 2017.
- [26] W. Chen, D. Byrd, S. Narayanan, and K. S. Nayak, "Intermittently tagged real-time MRI reveals internal tongue motion during speech production," *Magnetic Resonance in Medicine*, vol. 82, no. 2, pp. 600–613, 2019.
- [27] W. T. Fitch, "Vocal tract length and formant frequency dispersion correlate with body size in rhesus macaques," *The Journal of the Acoustical Society of America*, vol. 102, no. 2, pp. 1213–1222, 1997.
- [28] W. J. M. Peeters, "Diphthong dynamics: A cross-linguistic perceptual analysis of temporal patterns in Dutch, English, and German," 1996.
- [29] W. J. M. Peeters and W. J. Barry, "Diphthong dynamics: production and perception in Southern British English," in *EUROSPEECH*, 1989, pp. 1055–1058.