

Towards Accurate Phonetic Error Detection Through Phoneme Similarity Modeling

Xuanru Zhou¹, Jiachen Lian^{*2}, Cheol Jun Cho², Tejas Prabhune², Shuhe Li¹, William Li², Rodrigo Ortiz², Zoe Ezzes³, Jet Vonk³, Brittany Morin³, Rian Bogley³, Lisa Wauters³, Zachary Miller³, Maria Gorno-Tempini³, Gopala Anumanchipalli²

¹Zhejiang University, China ²UC Berkeley, United States

³UCSF, United States

^{*} Project Lead, corresponding to: jiachenlian@berkeley.edu

Abstract

Phonetic error detection, a core subtask of automatic pronunciation assessment, identifies pronunciation deviations at the phoneme level. Speech variability from accents and dysfluencies challenges accurate phoneme recognition, with current models failing to capture these discrepancies effectively. We propose a verbatim phoneme recognition framework using multi-task training with novel phoneme similarity modeling that transcribes what speakers actually say rather than what they're supposed to say. We develop and open-source *VCTK-accent*, a simulated dataset containing phonetic errors, and propose two novel metrics for assessing pronunciation differences. Our work establishes a new benchmark for phonetic error detection.

Index Terms: phonetic error detection, allophony, speech recognition, phoneme similarity modeling

1. Introduction

Speech pronunciation assessment plays a crucial role in language learning [1, 2] and diagnosis of speech disorders [3]. As traditional human assessment is time-consuming and lacks unified standards, recent advancement has shifted to Computer-Aided Pronunciation Training (CAPT) [4]. Parallel to the bloom of textual large language models [5], recent speech research focuses on audio-language foundation models [6, 7] that unify multiple tasks including pronunciation assessment. However, little evidence shows that multi-task learning improves performance in specific tasks like pronunciation modeling. Thus, *automatic pronunciation assessment remains a task-specific area* requiring domain-specific priors in model design.

Automatic pronunciation assessment involves detecting multiple aspects of speech quality [4], including fluency, intonation, phonetic errors, prosodic features, stress patterns, and rhythmic consistency. We focus on a core module: phonetic error detection, which aims to identify pronunciation deviations at a fine-grained phoneme level. This is fundamentally an *allophony* problem [8]. Specifically, if a person pronounces phoneme A, how confidently can we determine its relation to phoneme \hat{A} ? For instance, can AI models reliably distinguish between your pronunciation of the phoneme "TH" when the ground truth is "S" in the word *think*, as illustrated in Figure 1? So far, we have found no evidence of such capabilities in mainstream language learning platforms such as Duolingo [9], Speak [10], or even GPT-4o Voice [11].

Earlier studies have treated pronunciation allophony as a phoneme recognition or classification task, using normalized classification logits scores as confidence measures [3, 12, 13, 14, 15]. These are also called phonetic replacement errors in dysfluency modeling [13, 14, 16, 17, 18, 3, 15, 19, 20, 21, 22]. While [23] explored phoneme similarity in CTC training for

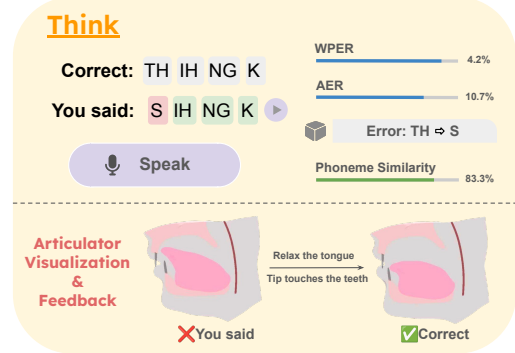


Figure 1: *Demo of our method. Model transcribes user speech into phoneme sequences, detects errors, scores using metrics, and generates articulator visualizations and feedback [25]: relax the tongue and place it against the roof of the mouth, with the tip lightly touching the teeth.*

dysarthric speech, their approach focuses on transcribing what speakers should have said rather than what they actually said, which is our research focus. Another relevant work [24] demonstrated that self-supervised speech representations can implicitly leverage phoneme similarity in a way that aligns more closely with human perception. However, it primarily focused on analysis and still fails to transcribe what the person actually said. Our work, on the other hand, focuses on the latter task, which is *distinct* and more *challenging*.

In this paper, we propose a framework for verbatim phoneme recognition with phoneme similarity modeling, trained using multi-task learning [26]: phoneme mapping and connectionist temporal classification (CTC), aiming to accurately transcribe the phonemes pronounced by the speaker. Phoneme similarity modeling, which serves as a method to reveal allophony information, better aligns with both how humans produce speech and how the human ear perceives it. We thus propose three novel phoneme similarity modeling methods: heuristic-based, articulatory-based, and Sylber-based [27] methods. Through this modeling, we compute the similarity score between each pair of phonemes. The process of integrating phoneme similarity into the training procedure is referred to as *soft-training*. To providing training material that incorporates phonetic errors and accurate labels, we follow the TTS-based simulation approach [16] and generate a simulated dataset with vowel and consonant substitutions, which we named *VCTK-accent* and open-sourced. In addition, We introduce two novel metrics: *WPER* and *AER*, which can be used not only for evaluating phoneme recognition models but also as pronunciation

assessment scores for speech.

Evaluated on the VCTK-accent test sets and real speech datasets, our method shows impressive performance. Our training strategy *multi-task learning & soft-training* significantly reduced the PER, WPER, and AER, proving its effectiveness. This highlights the crucial role of phoneme similarity modeling as a key approach to tackling the problem. Our work establishes a new benchmark for phonetic error detection. The project page is available at <https://berkeley-speech-group.github.io/Phonetic-Error-Detection/>.

2. Method

2.1. Data Simulation

To train the phoneme recognition model, accurate labels are essential, mapping directly to the word being pronounced. Thus, we follow the TTS-based dysfluency simulation pipeline described in [16]. First, we inject phonetic substitutions into each word from the text of the VCTK corpus [28] based on a predefined set of common phoneme substitution pairs listed in Table 1, which include vowel-to-vowel and consonant-to-consonant substitutions respectively. Next, we input the modified IPA sequences into the VITS [29] model to generate speech with phonetic errors. These resulting (*speech, modified phoneme sequence, reference word*) pairs are used as training data for our phoneme recognition model.

Table 1: Common CMU phoneme substitution pairs

Vowel			Consonant			
(AA, IY)	(AE, UW)	(AA, IH)	(P, G)	(T, ZH)	(K, B)	(M, S)
(OW, EH)	(AO, EH)	(UH, ER)	(N, SH)	(NG, F)	(L, T)	(R, D)
(AH, IY)	(ER, OW)	(AH, AE)	(W, K)	(TH, V)	(DH, Z)	(SH, HH)

2.2. Phoneme Similarity Modeling

Unlike [23] directly obtained the phoneme distance using PanPhon tool [30], we propose three methods for modeling phoneme similarity: heuristic-based, articulatory-based and Sylber-based methods. By incorporating perspectives from acoustics and syllables, these approaches more closely align with human pronunciation and auditory perception. We obtain a phoneme similarity matrix $S \in \mathbb{R}^{N \times N}$, where each value in the range (0, 1), indicating the similarity between each pair of phonemes. N represents the size of the phoneme dictionary.

2.2.1. Heuristic-based

Heuristic-based method calculates similarity by comparing phonemes based on eight features [31]: vowel or consonant, vowel length, vowel height, vowel frontness, lip rounding, consonant type, place of articulation, and consonant voicing. Each feature is assigned a normalized weight, and the similarity score between pairs of phonemes is computed by summing the weights of matching feature values. In this study, the weights are set as follows: 0.2, 0.1, 0.15, 0.15, 0.1, 0.2, 0.2, and 0.1, respectively. The visualization of the phoneme similarity matrix is presented in Figure 2.

2.2.2. Articulatory-based

We first construct a reference articulatory position for each phoneme using the VCTK corpus and an acoustic-to-articulatory inversion (AAI) model [25], employing a data-driven approach. Next, we compute the L2 distance between

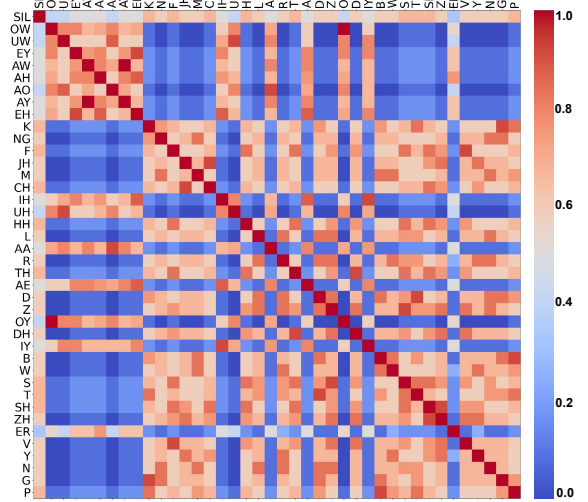


Figure 2: Heuristic-based phoneme similarity matrix

each pair of phonemes and apply min-max normalization to obtain the similarity scores.

2.2.3. Sylber-based

Human speech segmentation is naturally syllabic [27]. We then utilize the Sylber [27] feature, as it offers a clean and robust syllabic structure. First, we fine-tune Sylber using a phoneme classification task with a single linear classifier layer, employing the VCTK corpus. Then, following a similar approach outlined in Sec. 2.2.2, we construct a reference feature for each phoneme and compute the similarity score.

Overall, the heuristic-based method stems from the perspective of phoneme classification and definition. Both the articulatory and Sylber-based methods are data-driven approaches: the former emphasizes the acoustic aspect, while the latter focuses more on the syllabic aspect.

2.3. Verbatim Phoneme Recognition

We adopt the speech feature extracted from WavLM [32] to train a verbatim phoneme recognition model with multi-task learning and soft-training. The entire pipeline is illustrated in Figure 3, and the model architecture and training objectives are described in the following sections.

2.3.1. Model

The phoneme recognition model consists of a conformer [33] encoder, an autoregressive projection layer, and a linear classifier. This conformer encoder consists of 8 layers with 4 attention heads per layer. For the auto-regressive mechanism, For the autoregressive mechanism, at each timestep, the model combines the output features from the conformer encoder with the embedding of the previously predicted phoneme to predict the next phoneme, thereby effectively capturing the sequential dependencies in the data.

2.3.2. Multi-task Learning and Soft-training

We employ weighted loss-based multi-task training pipeline with two separate loss values: soft-CTC loss and soft-mapping loss. For the phoneme similarity matrix, each column(row)

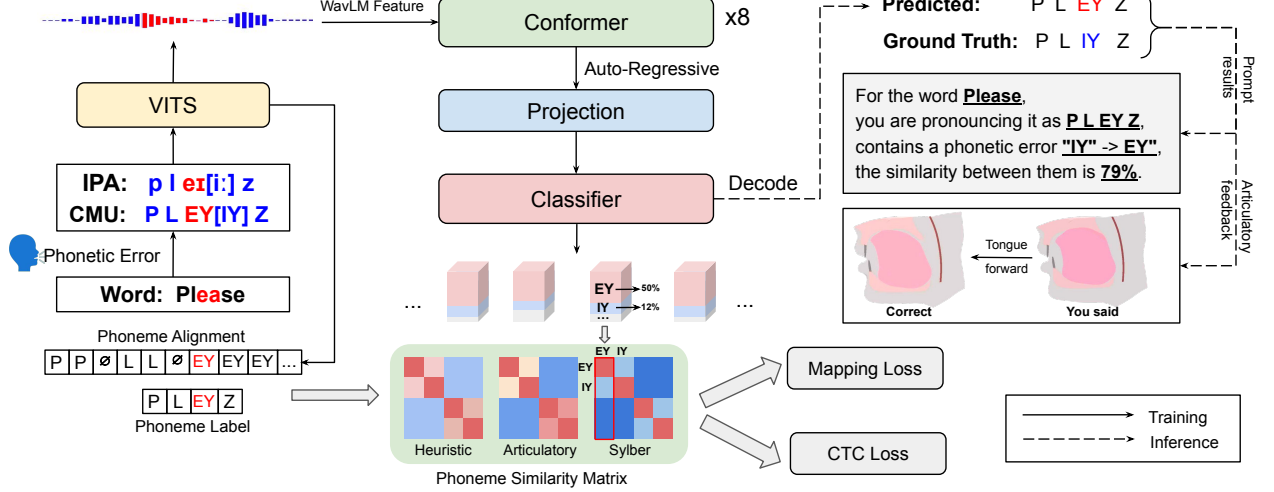


Figure 3: Pipeline of phoneme recognition and error detection: phonetic error of "IY" -> "EY" with a similarity score 79% in the word "Please", and the articulatory feedback is moving the tongue towards the front of the mouth.

represents the similarity score of this phoneme and all other phonemes. We can treat this vector as a *soft label* of this phoneme. For soft-CTC loss, we utilize the target phoneme's soft label to weight the emission probability, thereby reducing the penalty for prediction errors between similar phonemes. Traditional cross-entropy loss treats all phonemes as independent and ignores their similarities. Thus, in soft-mapping loss, we replace the one-hot encoded target with the target phoneme's soft label. The complete loss is shown below:

$$\begin{aligned} \mathbb{L} &= \lambda_{ctc} \cdot L_{sCTC} + \lambda_{map} \cdot L_{smap} \\ &= -\lambda_{ctc} \cdot \log \left(\sum_{z \in \mathcal{B}(y)} \prod_{t=1}^T \sum_{j=1}^N \hat{S}(z_t, j) \cdot p_t(j) \right) \\ &\quad + \lambda_{map} \cdot \sum_{t=1}^T \sum_{j=1}^N (p_t(j) - \hat{S}(y_t, j))^2 \end{aligned} \quad (1)$$

Where S denotes the normalized phoneme similarity matrix, $\mathcal{B}(y)$ denotes the set of all compatible alignments, y is the target phoneme sequence, p is model's output emission probability, and N is the phoneme dictionary size. In this work, we set $\lambda_{CTC} = 0.8$, $\lambda_{map} = 0.2$.

3. Experiment

3.1. Datasets

- **VCTK-Accent** is a TTS-based [16] simulated datasets, extended from VCTK corpus [28], which contains vowel and consonant phonetic errors, with simulation details provided in Sec. 2.1. The total duration of the dataset is 323.9 hours.
- **L2-ARCTIC** [34] includes recordings from 24 non-native English speakers, each recording about one hour of read speech from CMU's ARCTIC prompts. It provided forced-aligned phonetic transcriptions, and annotated 150 utterances per speaker for mispronunciation errors.
- **Speechocean762** [35] is an open-source non-native English speech corpus, which consists of 5000 English sentences. All the speakers are non-native, and their mother tongue is Mandarin, and half of the speakers are children.

- **MultiPA** [36] was collected from real-world open-response scenarios and consists of 50 audio clips, each lasting 10 to 20 seconds, from around 20 anonymous Mandarin-speaking users practicing English with a dialog chatbot.
- **PPA Speech** is collected from our clinical collaborators, and consists of recordings from 38 participants diagnosed with Primary Progressive Aphasia (PPA) [37]. We selected segments containing phonetic errors for evaluation.

For the above real speech datasets, we did the segmentation and annotation ourselves, the process is: we first segment out the words with phonetic error, then label the phoneme sequence that the speaker actually pronounced as the target.

3.2. Training

The phoneme recognition model is trained with 90/10 train/test split on VCTK-accent, with a batch size of 256. We use Adam optimization and gradient clipping, and the learning rate is $3e-4$. The model is trained for 30 epochs with total of 75 hours on an RTX A6000.

3.3. Evaluation Metrics

3.3.1. Phoneme Error Rate (PER)

PER measure of how many errors (inserted, deleted, and substitute phonemes) are predicting phoneme sequences compared to the actual phoneme sequence. It calculated by dividing the number of phoneme errors by the total number of phonemes.

3.3.2. Weight Phoneme Error Rate (WPER)

The single-word phoneme error rate (PER) is a relatively coarse metric, as it only reflects the presence or absence of phonetic errors, lacking the ability to capture nuanced pronunciation differences. To address this limitation, we introduce a more refined metric: the Weighted Phoneme Error Rate (WPER). In the case of substitutions, we replace the count of substitutions with the sum of the phoneme similarities between the substituted pairs. The equation is shown below:

$$WPER = \frac{D + \sum_{(p_r, p_s)} (1 - S(p_r, p_s)) + I}{L} \quad (2)$$

Table 2: Evaluation on VCTK-accent testsets with three different phoneme similarity modeling methods

Metrics	Vowel			Consonant			All		
	Heuristic	Articulatory	Sylber	Heuristic	Articulatory	Sylber	Heuristic	Articulatory	Sylber
PER (% , ↓)	12.39	10.15	10.04	13.93	13.62	12.89	15.85	12.37	16.51
WPER (% , ↓)	8.19	7.34	8.83	9.18	7.75	10.42	9.39	7.41	9.57
AER (% , ↓)	9.22	7.93	9.37	10.84	10.76	11.09	12.82	10.53	12.08

Where L be the length of the reference phoneme sequence, D and I the counts of deletions and insertions, S the phoneme similarity matrix constructed in Sec. 2.2, p_r , p_s are the reference and substitute phoneme, respectively.

3.3.3. Articulatory Error Rate (AER)

We also propose a metric that considers the articulatory distance between different phonemes. For each speech sample, we apply the AAI model [25] to convert the waveform into articulatory features. We then calculate the L2 distance between the articulatory features of the current frame and the target phoneme, using the mapping we constructed in Sec. 2.2.2. If the distance exceeds a threshold, denoted as τ , the frame is classified as negative. AER is then computed as the ratio of negative frames to the total length of the speech. In this work, we set τ to be 0.5 times the distance between the two most distant phonemes.

Table 3: Phoneme recognition with different training tasks

Training tasks	PER (% , ↓)	WPER (% , ↓)	AER (% , ↓)
w/o soft-training	18.31	16.98	17.42
+ smap	14.81	11.98	11.87
+ sCTC	16.14	14.36	13.29
+ smap + sCTC	12.37	7.41	10.53

Table 4: Evaluation on real speech datasets

Datasets	PER (% , ↓)	WPER (% , ↓)	AER (% , ↓)
L2-ARCTIC [34]	26.67	16.53	18.79
Speechocean762 [35]	28.33	17.74	16.33
MultiPA [36]	30.49	19.01	20.45
PPA Speech [37]	38.63	21.78	21.45

3.4. Validation

To assess the performance of the phoneme recognition model, we first conduct evaluation on the VCTK-accent dataset, using three different phoneme similarity modeling methods, and focusing on vowel, consonant, and whole testsets, respectively. The results are presented in Table 2. Furthermore, to comprehensively evaluate the impact of phoneme similarity modeling, we trained three additional models with different training objectives: multi-task learning without soft-training, adding soft phoneme mapping training, and adding soft CTC training. For the latter two, the soft training used the articulatory-based method (the best, as referenced in Table 2). We then compared these results with the benchmark (smap + sCTC), as shown in Table 3. Additionally, to assess model’s generalizability in real-world scenarios, we validated the model on L2-Arctic, Speechocean762, MultiPA, and PPA speech, with the results presented in Table 4.

3.5. Results and Discussion

Since there are no established benchmarks suitable for direct comparison, we report PER along with our proposed WPER and AER for each evaluation. As shown in Table 2, the phoneme recognition model’s performance on the VCTK-accent testset yields PER consistently around 10–15%, indicating a solid performance level. Meanwhile, WPER and AER range between 7–12%, providing a more reasonable measure of the speech pronunciation differences. Notably, vowel substitutions appear to be transcribed more accurately than consonant substitutions, possibly due to the higher TTS simulation quality for vowel substitutions. Among the three phoneme similarity modeling methods, the articulatory-based approach demonstrates the best performance, which aligns with the fundamental principles of human speech production.

From Table 3, we can observe that multi-task learning with soft training yields the best results, as it outperforms the other tasks across all three metrics. It reduces the PER from 18.31% to 12.37%, and the WPER and AER decrease by approximately 8%. Among the individual tasks, only adding soft phoneme mapping and soft CTC training show improvements compared to the task without soft training, with soft phoneme mapping demonstrating a greater improvement. As indicated in Table 4, the model’s performance at real speech datasets are not as strong compared to the VCTK-accent test set, but the generalization still at a good level. Among the datasets, the model performs best on L2-Arctic and Speechocean762, which contain read speech texts. Its performance is relatively lower on MultiPA, which features open-response spontaneous speech, and is weakest on PPA speech, representing the most severe form of pronunciation error. The performance aligns with the difficulty level of the speech.

4. Limitation and Conclusion

This paper presents a framework for verbatim phoneme recognition and error detection using multi-task learning and soft training, incorporating phoneme similarity modeling. Results show that modeling similarity significantly improves transcription accuracy and sets a new benchmark for phonetic error detection. Nonetheless, limitations remain. CMU phonemes may not align well with actual speech; future work could explore IPA [38] or syllables. TTS quality also needs refinement, and LLMs could generate more natural substitution pairs. Finally, integrating articulatory feedback in gestural [3, 39, 40, 41], robust visual units [42] and kinematic spaces may enable more human-aligned, closed-loop pronunciation learning.

5. Acknowledgements

Thanks for support from UC Noyce Initiative, Society of Hellman Fellows, NIH/NIDCD, and the Schwab Innovation fund.

6. References

- [1] Y. El Kheir, A. Ali, and S. A. Chowdhury, "Automatic pronunciation assessment - a review," in *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 8304–8324.
- [2] K. Truong, A. Neri, C. Cucchiari, and H. Strik, "Automatic pronunciation error detection: an acoustic-phonetic approach," in *InSTIL/ICALL 2004 Symposium on Computer Assisted Learning*, 2004, p. paper 032.
- [3] J. Lian, X. Zhou, Z. Ezzes, J. Vonk, B. Morin, D. P. Baquirin, Z. Miller, M. L. Gorno Tempini, and G. Anumanchipalli, "Ssdm: Scalable speech dysfluency modeling," in *Advances in Neural Information Processing Systems*, vol. 37, 2024.
- [4] A. Lee *et al.*, "Language-independent methods for computer-assisted pronunciation training," Ph.D. dissertation, Massachusetts Institute of Technology, 2016.
- [5] OpenAI, "Chatgpt," 2022. [Online]. Available: <https://openai.com/chatgpt/>
- [6] S. Ji, Y. Chen, M. Fang, J. Zuo, J. Lu, H. Wang, Z. Jiang, L. Zhou, S. Liu, X. Cheng *et al.*, "Wavchat: A survey of spoken dialogue models," *arXiv preprint arXiv:2411.13577*, 2024.
- [7] C.-y. Huang, W.-C. Chen, S.-w. Yang, A. T. Liu, C.-A. Li, Y.-X. Lin, W.-C. Tseng, A. Diwan, Y.-J. Shih, J. Shi *et al.*, "Dynamic-superb phase-2: A collaboratively expanding benchmark for measuring the capabilities of spoken language models with 180 tasks," *ICLR*, 2025.
- [8] A. Boomershire, K. C. Hall, E. Hume, and K. Johnson, "The impact of allophony versus contrast on speech perception," *Contrast in phonology*, pp. 143–172, 2008.
- [9] Duolingo, "Duolingo: Language learning app," 2025. [Online]. Available: <https://www.duolingo.com/>
- [10] Speak, "Speak: Ai-powered language learning," 2025. [Online]. Available: <https://www.speak.com/>
- [11] OpenAI, "Gpt-4o real time," 2024. [Online]. Available: <https://openai.com/index/hello-gpt-4o/>
- [12] M. Yang, K. Hirschi, S. D. Looney, O. Kang, and J. H. Hansen, "Improving mispronunciation detection with wav2vec2-based momentum pseudo-labeling for accentedness and intelligibility assessment," in *Interspeech 2022*, 2022, pp. 4481–4485.
- [13] J. Lian, C. Feng, N. Farooqi, S. Li, A. Kashyap, C. J. Cho, P. Wu, R. Netzorg, T. Li, and G. K. Anumanchipalli, "Unconstrained dysfluency modeling for dysfluent speech transcription and detection," in *2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, 2023, pp. 1–8.
- [14] J. Lian and G. Anumanchipalli, "Towards hierarchical spoken language disfluency modeling," in *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics*, Mar. 2024, pp. 539–551.
- [15] J. Lian, X. Zhou, Z. Ezzes, J. Vonk, B. Morin, D. Baquirin, Z. Mille, M. L. G. Tempini, and G. K. Anumanchipalli, "Ssdm 2.0: Time-accurate speech rich transcription with non-fluencies," *arXiv preprint arXiv:2412.00265*, 2024.
- [16] X. Zhou, A. Kashyap, S. Li, A. Sharma, B. Morin, D. Baquirin, J. Vonk, Z. Ezzes, Z. Miller, M. Tempini, J. Lian, and G. Anumanchipalli, "Yolo-stutter: End-to-end region-wise speech dysfluency detection," in *Interspeech 2024*, 2024, pp. 937–941.
- [17] X. Zhou, C. J. Cho, A. Sharma, B. Morin, D. Baquirin, J. Vonk, Z. Ezzes, Z. Miller, B. L. Tee, M. L. Gorno-Tempini *et al.*, "Stutter-solver: End-to-end multi-lingual dysfluency detection," in *2024 IEEE Spoken Language Technology Workshop (SLT)*. IEEE, 2024, pp. 1039–1046.
- [18] X. Zhou, J. Lian, C. J. Cho, J. Liu, Z. Ye, J. Zhang, B. Morin, D. Baquirin, J. Vonk, Z. Ezzes, Z. Miller, M. L. G. Tempini, and G. Anumanchipalli, "Time and tokens: Benchmarking end-to-end speech dysfluency detection," 2024. [Online]. Available: <https://arxiv.org/abs/2409.13582>
- [19] J. Zhang, X. Zhou, J. Lian, S. Li, W. Li, Z. Ezzes, R. Bogley, L. Wauters, Z. Miller, J. Vonk, B. Morin, M. Gorno-Tempini, and G. Anumanchipalli, "Analysis and evaluation of synthetic data generation in speech dysfluency detection," *Interspeech*, 2025.
- [20] C. Guo, J. Lian, X. Zhou, J. Zhang, S. Li, Z. Ye, H. J. Park, A. Das, Z. Ezzes, J. Vonk, B. Morin, R. Bogley, L. Wauters, Z. Miller, M. Gorno-Tempini, and G. Anumanchipalli, "Dysfluent wfst: A framework for zero-shot speech dysfluency transcription and detection," *Interspeech*, 2025.
- [21] Z. Ye, J. Lian, X. Zhou, J. Zhang, H. Li, S. Li, C. Guo, A. Das, P. Park, Z. Ezzes, J. Vonk, B. Morin, R. Bogley, L. Wauters, Z. Miller, M. Gorno-Tempini, and G. Anumanchipalli, "Seamless dysfluent speech text alignment for disordered speech analysis," *Interspeech*, 2025.
- [22] J. Lian, X. Zhou, C. Guo, Z. Ye, Z. Ezzes, J. M. Vonk, B. Morin, D. Baquirin, Z. Miller, M. L. Gorno-Tempini, and G. K. Anumanchipalli, "Automatic detection of articulatory-based disfluencies in primary progressive aphasia," *IEEE Journal of Selected Topics in Signal Processing*, pp. 1–17, 2025.
- [23] W. Lee, S. Im, H. Do, Y. Kim, J. Ok, and G. G. Lee, "Dypcl: Dynamic phoneme-level contrastive learning for dysarthric speech recognition," *NAACL*, 2025.
- [24] K. Choi, E. Yeo, K. Chang, S. Watanabe, and D. Mortensen, "Leveraging allophony in self-supervised speech models for atypical pronunciation assessment," in *NAACL*, 2025.
- [25] C. J. Cho, P. Wu, T. S. Prabhune, D. Agarwal, and G. K. Anumanchipalli, "Coding speech through vocal tract kinematics," in *IEEE JSTSP*, 2025.
- [26] M. Crawshaw, "Multi-task learning with deep neural networks: A survey," 2020. [Online]. Available: <https://arxiv.org/abs/2009.09796>
- [27] C. J. Cho, N. Lee, A. Gupta, D. Agarwal, E. Chen, A. W. Black, and G. K. Anumanchipalli, "Sylber: Syllabic embedding representation of speech from raw audio," in *ICLR*, 2025.
- [28] J. Yamagishi, C. Veaux, K. MacDonald *et al.*, "Cstr vctk corpus: English multi-speaker corpus for cstr voice cloning toolkit (version 0.92)," *University of Edinburgh. The Centre for Speech Technology Research (CSTR)*, 2019.
- [29] J. Kim, J. Kong, and J. Son, "Conditional variational autoencoder with adversarial learning for end-to-end text-to-speech," in *International Conference on Machine Learning*, 2021.
- [30] D. R. Mortensen, P. Littell, A. Bharadwaj, K. Goyal, C. Dyer, and L. Levin, "PanPhon: A resource for mapping IPA segments to articulatory feature vectors," in *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*. Osaka, Japan: The COLING 2016 Organizing Committee, Dec. 2016, pp. 3475–3484.
- [31] C. Anderson, *Essentials of linguistics*. McMaster University, 2018.
- [32] S. Chen, C. Wang, Z. Chen, Y. Wu, S. Liu, Z. Chen, J. Li, N. Kanda, T. Yoshioka, X. Xiao, J. Wu, L. Zhou, S. Ren, Y. Qian, Y. Qian, J. Wu, M. Zeng, X. Yu, and F. Wei, "Wavlm: Large-scale self-supervised pre-training for full stack speech processing," *IEEE JSTSP*, 2022.
- [33] A. Gulati, J. Qin, C.-C. Chiu, N. Parmar, Y. Zhang, J. Yu, W. Han, S. Wang, Z. Zhang, Y. Wu, and R. Pang, "Conformer: Convolution-augmented transformer for speech recognition," in *Interspeech 2020*, 2020, pp. 5036–5040.
- [34] G. Zhao, S. Sonsaat, A. Silpachai, I. Lucic, E. Chukharev-Hudilainen, J. Levis, and R. Gutierrez-Osuna, "L2-arctic: A non-native english speech corpus," in *Interspeech*, 2018.
- [35] J. Zhang, Z. Zhang, Y. Wang, Z. Yan, Q. Song, Y. Huang, K. Li, D. Povey, and Y. Wang, "speechoccean762: An open-source non-native english speech corpus for pronunciation assessment," in *Interspeech 2021*, 2021, pp. 3710–3714.
- [36] Y.-W. Chen, Z. Yu, and J. Hirschberg, "Multipa: A multi-task speech pronunciation assessment model for open response scenarios," in *Interspeech 2024*, 2024, pp. 297–301.

- [37] M. L. Gorno-Tempini, A. E. Hillis, S. Weintraub, A. Kertesz, M. Mendez, S. F. Cappa, J. M. Ogar, J. D. Rohrer, S. Black, B. F. Boeve *et al.*, “Classification of primary progressive aphasia and its variants,” *Neurology*, vol. 76, no. 11, pp. 1006–1014, 2011.
- [38] Wikipedia contributors, “International phonetic alphabet.” [Online]. Available: https://en.wikipedia.org/wiki/International_Phonetic_Alphabet
- [39] J. Lian, A. W. Black, L. Goldstein, and G. K. Anumanchipalli, “Deep Neural Convolutional Matrix Factorization for Articulatory Representation Decomposition,” in *Proc. Interspeech 2022*, 2022, pp. 4686–4690.
- [40] J. Lian, A. W. Black, Y. Lu, L. Goldstein, S. Watanabe, and G. K. Anumanchipalli, “Articulatory representation learning via joint factor analysis and neural matrix factorization,” in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023, pp. 1–5.
- [41] P. Wu, T. Li, Y. Lu, Y. Zhang, J. Lian, A. W. Black, L. Goldstein, S. Watanabe, and G. K. Anumanchipalli, “Deep Speech Synthesis from MRI-Based Articulatory Representations,” in *Proc. INTERSPEECH 2023*, 2023, pp. 5132–5136.
- [42] J. Lian, A. Baevski, W.-N. Hsu, and M. Auli, “Av-data2vec: Self-supervised learning of audio-visual speech representations with contextualized target representations,” in *2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2023, pp. 1–8.