

TOWARDS FAIR ASR FOR SECOND LANGUAGE SPEAKERS USING FAIRNESS PROMPTED FINETUNING

Monorama Swain¹, Bubai Maji², Jagabandhu Mishra³, Markus Schedl¹,
Anders Søgaard⁴, Jesper Rindom Jensen⁵

¹Johannes Kepler University Linz, Austria, ²IIT Kharagpur, India, ³University of Eastern Finland,
⁴University of Copenhagen, Denmark, ⁵Aalborg University, Denmark

ABSTRACT

In this work, we address the challenge of building fair English ASR systems for second-language speakers. Our analysis of widely used ASR models, Whisper and Seamless-M4T, reveals large fluctuations in word error rate (WER) across 26 accent groups, indicating significant fairness gaps. To mitigate this, we propose fairness-prompted finetuning with lightweight adapters, incorporating Spectral Decoupling (SD), Group Distributionally Robust Optimization (Group-DRO), and Invariant Risk Minimization (IRM). Our proposed fusion of traditional empirical risk minimization (ERM) with cross-entropy and fairness-driven objectives (SD, Group DRO, and IRM) enhances fairness across accent groups while maintaining overall recognition accuracy. In terms of macro-averaged word error rate, our approach achieves a relative improvement of 58.7% and 58.5% over the large pretrained Whisper and Seamless-M4T, and 9.7% and 7.8% over them, finetuning with standard empirical risk minimization with cross-entropy loss.

Index Terms— Automatic speech recognition, fairness in speech recognition, accent and language variation

1. INTRODUCTION

Speech technologies can be a driver of equality, making information accessible to social groups with limited access: dyslexics, non-literates, congenitally blind, children, the elderly, and language learners, if they work well across all social groups, including non-native speakers [1]. Performance drops for non-native speakers can often be traced back to the influence of the speaker’s mother tongue and its effect on pronunciation [2], [3]. Accents may exhibit a lot of variation: In Indian English, for example, the /z/ sound is often replaced by /s/, especially in the speech of those whose native languages do not distinguish between the two. Moreover, *very* might be pronounced as [ˈve:ri], using the retroflex /ó/ instead of the alveolar /r/ [4].

Automatic speech recognition (ASR) systems have achieved remarkable performance in high-resource languages such as English. Yet, several studies have shown that these systems exhibit unequal performance across demographic and accent groups [5, 6, 7, 8, 9], a disparity further illustrated

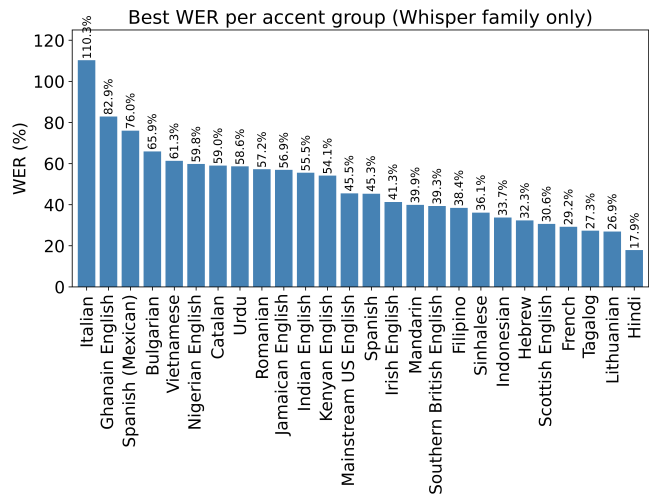


Fig. 1. Best WER (%) per accent group for Whisper models before fine-tuning. For each group, the WER shown corresponds to the lowest WER obtained among Whisper variants (i.e., tiny, base, small, medium, large)

in Figure 1. Such fairness gaps highlight the pressing need for systematic methodologies to evaluate and mitigate bias, aligning with the call in [9] that “Every voice matters.” A key reason for these disparities is the over-representation of mainstream accents in training and evaluation data [10], leading to substantial performance gaps for underrepresented accents.

In this paper, we focus on improving the fairness and robustness of ASR across various English accents, with particular attention to *second-language* (L2) speakers. There has been some attempts already, however, to study fairness in ASR. For example, in [11], a mixed-effects poisson regression has been introduced as a statistical approach to measure and interpret word error rate (WER) differences among subgroups. This method has been shown to address challenges such as handling unobserved heterogeneity across speakers and identifying the sources of WER gaps between subgroups. Similarly, [12] introduced counterfactual training methods for Connectionist Temporal Classification-based ASR, where demographic information such as gender, age, education are counterfactually

modified while the linguistic content remains fixed. Their method helps to reduce the variance in character error rate across demographic groups.

Recent work, such as [13], has benchmarked ASR performance across languages and accents, highlighting performance disparities of foundational ASR models across accent groups. Similarly, [5] leveraged one-hot accent embeddings to improve adaptation to diverse speakers. However, prior studies have primarily focused on modeling accent variation, and to the best of our knowledge, none have explored *fairness prompting* as a means to fine-tune foundational models. Our work addresses this gap by explicitly targeting L2 English varieties and incorporating *fairness disparities* into the learning objective, thereby enabling *fair finetuning*. Our main contributions are:

- We conducted systematic fairness analyzes targeting L2 English speakers, highlighting disparities that remain hidden when treating all accents uniformly.
- We evaluated two families of off-the-shelf speech models (Whisper and SeamlessM4T) under standard vanilla fine-tuning (i.e., empirical risk minimization (ERM)) as the baseline with the fairness-promoting algorithms: Group-DRO; spectral decoupling (SD); and invariant risk minimization (IRM).
- Finally, we propose a novel fusion of ERM, SD, Group-DRO, and IRM (see Figure 2) to evaluate its effectiveness in reducing performance gaps across 26 L2 English accents. In addition, we also analyze how model size and linguistic/data properties influence the effectiveness of fairness-promoting fine-tuning.

2. METHODOLOGY

2.1. Problem Statement

Our goal is to minimize the disparities in ASR performance across 26 groups. Each group contains English utterances shaped by a distinct L2 corresponding to one native language. Thus, each sample consists of the information (x, y, g) , where $x \in \mathbb{R}^T$ is the speech sample, where T is the length of the input speech signal, $y = (y_1, \dots, y_M)$ is the English transcription sequence, and $g \in \{1, \dots, 26\}$ denotes the accent group.

2.2. ASR Model

We adopt large-scale pretrained sequence-to-sequence ASR models, i.e., Whisper family and SeamlessM4T:

Whisper [14]: A family set of five ASR models (tiny, base, small, medium, and large)¹, which are trained in the same way, but of different sizes, developed by OpenAI. The Whisper models are trained on 680,000 hours of data sourced from the web in 97 languages. During training, Whisper uses data

¹All Whisper model variants are available at <https://huggingface.co/models?search=openai/whisper>.

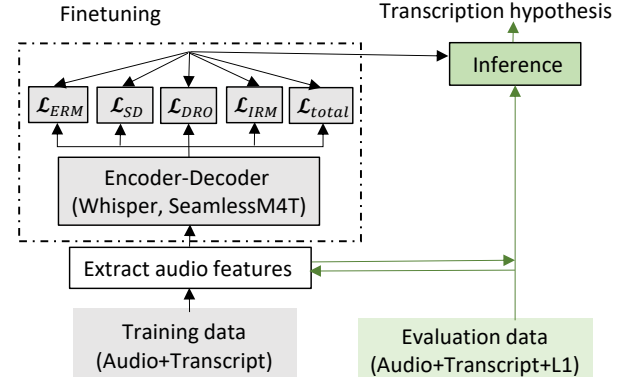


Fig. 2. The fairness-prompting speech recognition system.

augmentation techniques that transform audio spectrograms by applying methods like time warping, frequency masking, and time masking. These techniques enhance the model’s ability to generalize across various acoustic environments. Whisper is built as a traditional encoder-decoder system and is trained in a weakly supervised manner on unreleased multilingual audio data. To handle long audio, the system divides the audio into 30-second chunks for processing.

SeamlessM4T [15]: is a model developed by MetaAI that supports speech-to-speech translation, speech-to-text translation, text-to-speech translation, text-to-text translation, and ASR for up to 100 languages. The model has two versions²: Medium with 1.2 billion parameters and a size of 6.4 GB, and Large with 2.3 billion parameters and a size of 10.7 GB.

2.3. Finetuning with ERM

ERM [16]: The ERM baseline minimizes the empirical mean cross-entropy over the pooled training set, which helps to maximize the likelihood of the reference transcription. Accordingly, the ERM objective is formulated as:

$$\mathcal{L}_{\text{ERM}} = \frac{1}{N} \sum_{n=1}^N \mathcal{L}_{\text{ASR}}^{(n)} \quad (1)$$

where N is the total number of training utterances and $\mathcal{L}_{\text{ASR}}^{(n)}$ is the ASR loss (cross-entropy) for the n -th utterance. This standard training objective, which optimizes average loss, tends to favor majority groups and can result in higher WER for underrepresented accents.

2.4. Fairness-Promoting Finetuning

To adapt the considered models fairly across groups, we fine-tune with a multi-objective fairness loss. This allows adaptation to non-native English accents.

SD [17]: SD is a regularization method that penalizes the squared logit magnitude. It helps to reduce overconfident

²Both Seamless models are available at <https://huggingface.co/models?search=seamless>.

predictions and spurious correlations to improve generalization across accents. The SD loss is computed as:

$$\mathcal{L}_{\text{SD}} = \mathcal{L}_{\text{ERM}} + \lambda \|o\|_2^2 \quad (2)$$

where o are the model logits before the Softmax layer and λ is a regularization coefficient.

Group-DRO [18]: Standard ERM minimizes average risk, which can improve majority accents while ignoring minority groups. Group-DRO shifts optimization to the *worst-performing group*, explicitly improving fairness. This can be mathematically expressed as:

$$\mathcal{L}_{\text{DRO}} = \max_{g \in G} \mathcal{L}_{\text{ASR}}^{(g)} \quad (3)$$

where G is the set of accent groups and $\mathcal{L}_{\text{ASR}}^{(g)}$ is the average ASR loss within group g .

IRM [19]: IRM encourages the model to learn consistent predictors across environments (e.g., channel conditions or demographic metadata) to reduce spurious correlations. The IRM objective is defined as:

$$\mathcal{L}_{\text{IRM}} = \sum_{e \in E} \|\nabla_w \mathcal{L}_{\text{ASR}}^{(e)}(w \cdot \Phi)\|^2 \quad (4)$$

where E denotes the set of environments, $\mathcal{L}_{\text{ASR}}^{(e)}$ the ASR loss in environment e , Φ the learned features, w is a scalar classifier.

Fusion: The fusion objective reduces performance gaps between accent groups, yielding a fairer ASR system for English as L2. The overall objective is a weighted combination of all terms and is computed as:

$$\mathcal{L}_{\text{total}} = \lambda_e \mathcal{L}_{\text{ERM}} + \lambda_s \mathcal{L}_{\text{SD}} + \lambda_d \mathcal{L}_{\text{DRO}} + \lambda_i \mathcal{L}_{\text{IRM}} \quad (5)$$

where $\lambda_e, \lambda_s, \lambda_d, \lambda_i$ are scalar weights that balance the contributions of each loss.

3. EXPERIMENTAL SETUP AND RESULTS

3.1. Dataset

We used the Edinburgh International Accents of English Corpus (EdAcc) [20] – a new ASR dataset composed of 40 hours of English dyadic conversations between speakers with a diverse set of accents. EdAcc contains more than 40 self-reported English accents from speakers of 51 different first languages. It includes 26 distinct first- and second-language varieties of English, along with a linguistic background profile of each speaker, detailing how long they have spoken English and where they have lived. The conversations range from 20–60 minutes in duration. In our experiments, we used standard data splits as described in the dataset description [20].

3.2. Experimental Setup

For a fair comparison, we evaluated all four fine-tuning strategies, including our baseline ERM approach. We set

Table 1. Comparison of micro-average WER (left of “/”) and min–max difference (right of “/”) across Whisper and Seamless models under different fine-tuning strategies.

Model	w/o FT	ERM	DRO	SD	IRM	Fusion
Whisper						
Large	58.3 / 114.0	26.7 / 30.1	33.3 / 37.6	35.7 / 31.5	38.3 / 53.2	24.1 / 30.8
Medium	61.1 / 92.6	31.8 / 47.6	38.2 / 48.3	40.2 / 55.4	42.6 / 59.9	28.2 / 35.3
Small	67.8 / 120.5	32.9 / 39.4	38.3 / 41.9	41.4 / 44.3	43.8 / 53.2	30.3 / 45.1
Base	65.5 / 103.6	38.7 / 47.9	44.3 / 43.0	48.9 / 57.1	41.4 / 59.2	33.7 / 59.8
Tiny	96.0 / 124.0	48.9 / 54.1	51.4 / 56.9	57.8 / 64.0	55.7 / 85.1	41.0 / 57.8
Seamless						
Large	65.3 / 52.7	29.4 / 43.3	28.1 / 36.0	26.3 / 28.5	32.8 / 48.6	27.1 / 37.6
Medium	67.2 / 42.5	40.5 / 50.8	26.8 / 38.2	28.3 / 34.2	36.3 / 48.7	29.0 / 29.0

$\lambda_e = \lambda_d = 1$ to retain empirical risk and Group-DRO as the primary optimization, while $\lambda_s = 0.06$ and $\lambda_i = 0.01$ were selected through a greedy search in the range 0.01–1. The learning rate was fixed at 4×10^{-5} . These hyperparameters were tuned on validation performance to balance fairness improvements with ASR accuracy. The ERM baseline adapter was trained on top of the pre-trained Whisper model, enabling adaptation to the EdAcc dataset.

3.3. Fairness-Oriented Evaluation Metrics

To quantify fairness across accent groups, we report (i) the Word Error Rate for each accent group ($\text{WER}(g)$), (ii) the Macro-average of $\text{WER}(g)$ [21], showing the overall ASR accuracy, averaged fairly across groups (see Equation 6), and (iii) the Min-Max gap of $\text{WER}(g)$ [22], to evaluate the disparity between the worst and best performing groups (see Equation 7).

$$\text{Macro-average} = \frac{1}{|G|} \sum_{g \in G} \text{WER}(g) \quad (6)$$

$$\text{Min-Max Gap} = \max_{g \in G} \text{WER}(g) - \min_{g \in G} \text{WER}(g) \quad (7)$$

where G is the total number of accent groups (=26).

3.4. Results and Discussion

Table 1 reports the micro-average WER and min–max disparities for the Whisper and Seamless families under different fine-tuning strategies. Without fine-tuning, both families perform poorly on accented English. Whisper yields particularly high WERs (58–96%) and large disparities (92–124%), while Seamless shows smaller disparities (42–53%) at baseline, reflecting stronger multilingual pretraining. Fine-tuning consistently improves performance. ERM reduces WER by nearly 50% relative to the baseline—for instance, Whisper-large drops from 58.3% WER to 26.7%. Fairness-oriented methods such as DRO and SD further reduce disparities, with Seamless-large under SD achieving the smallest gap (28.5%).

Our proposed fusion objective outperforms the individual fairness-oriented objectives. Across all Whisper architectures, fusion lowers error rates across accents, achieving

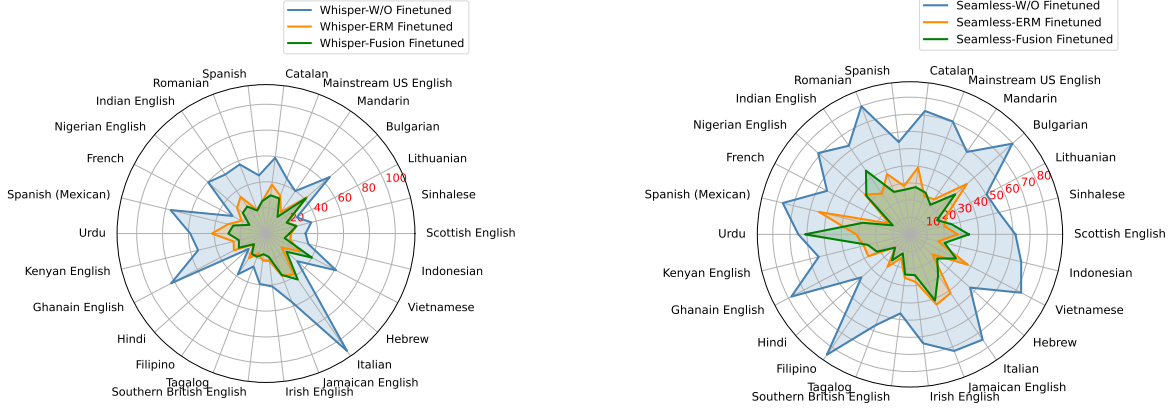


Fig. 3. WER (%) across accent groups for Whisper (left) and Seamless (right) models under different fine-tuning strategies. Without finetuning (blue) show high WER and disparities. ERM (orange) reduces errors but retains variability, while Fusion (green) yields the most balanced profiles, lowering error rates for challenging accents and improving fairness across groups.

the best performance with Whisper-large at 24.1% WER. In Seamless, fusion produces a more balanced profile; for example, Seamless-Medium fusion achieves 29% WER and a 29% min–max gap, compared to Seamless-large with 27.1% WER and a larger 37.6% gap. Thus, Whisper provides stronger performance in terms of macro-average WER, whereas Seamless maintains smaller min–max disparities across accents.

Figures 3 illustrate per-accent WER for Whisper and Seamless. Without fine-tuning, both models show large spikes for under-represented accents, consistent with the disparities in Table 1. ERM reduces error rates but leaves several accents with high WER. Our proposed fusion strategy improves performance across all accents for Whisper, and for Seamless, it improves almost all accents, with the exception of Urdu and Indian English. This limitation may stem from the under-representation of these accent groups in the initial pre-training, and we plan to investigate this further in future work.

3.5. Analysis

Effect of Model Size: The impact of fairness-promoting algorithms on performance decreases as models grow larger. Larger models reduce WER across all strategies, but the gap between ERM, DRO, SD, and Fusion narrows with scale. This observation, shown in Figure 4, suggests that scaling itself mitigates some fairness trade-offs while raising new research for sustainable language modeling.

Typological Distance: We consistently observed the highest performance in Asian accents, followed by European accents (e.g., Romanian), and African accents, with some more European varieties (e.g., Italian). We found no clear correlation with typological distance. For example, Romanian and Italian are closely related but differed by nearly 20 WER points.

Word Length: We also computed the average word length in the test data for each variety of English to control for this potential confounding factor. However, we found no correlation between word length and WER.

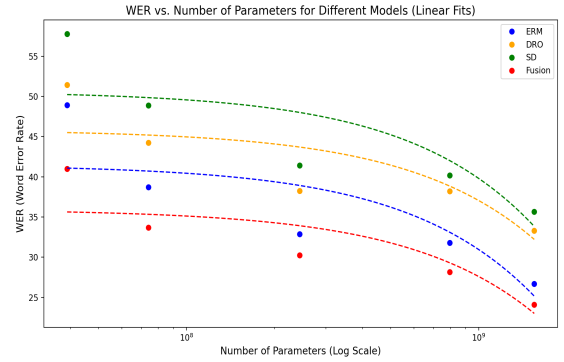


Fig. 4. WER across Whisper models sizes.

4. CONCLUSION

We investigated the performance gaps of widely adopted ASR models such as Whisper and Seamless on English accented by L2 speakers. To address these disparities, we applied fairness-promoting objectives—SD, DRO, and IRM—for fine-tuning. Our results show that fine-tuning with fairness objectives consistently improves macro-averaged WER compared to both pretrained models and standard empirical risk minimization. Moreover, our proposed fusion of ERM with fairness objectives further enhances performance, outperforming individual objectives. We also observed that increasing model size reduces micro-averaged WER, highlighting the benefit of scaling. Although we explored potential correlations between WER and factors such as native language typological distance and average word length, we did not find strong evidence of direct relationships. Nevertheless, our fine-tuning approach delivered consistent improvements across most L2 accents, with only a few challenging cases such as Urdu and Indian English in the case of Seamless. In future work, we will focus on these difficult accents to further reduce the performance gap and advance towards more equitable ASR systems.

5. REFERENCES

- [1] T.-P. Tan and L. Besacier, "Acoustic model interpolation for non-native speech recognition," in *2007 IEEE International Conference on Acoustics, Speech and Signal Processing-ICASSP'07*, vol. 4, 2007, pp. IV–1009.
- [2] K. Livescu and J. Glass, "Lexical modeling of non-native speech for automatic speech recognition," in *2000 IEEE international conference on acoustics, speech, and signal processing. Proceedings (Cat. No. 00CH37100)*, vol. 3. IEEE, 2000, pp. 1683–1686.
- [3] Y. R. Oh, J. S. Yoon, and H. K. Kim, "Acoustic model adaptation based on pronunciation variability analysis for non-native speech recognition," *Speech Communication*, vol. 49, no. 1, pp. 59–70, 2007.
- [4] R. Gargesh, "Indian english: Phonology," *Varieties of English*, vol. 4, no. 2, pp. 231–243, 2008.
- [5] P. DHERAM, M. Ramakrishnan, A. Raju, I.-F. Chen, B. King, K. Powell, M. Saboowala, K. Shetty, and A. Stolcke, "Toward fairness in speech recognition: Discovery and mitigation of performance disparities," in *Proc. Interspeech 2022*, 2022, pp. 1268–1272.
- [6] A. Koenecke, A. Nam, E. Lake, J. Nudell, M. Quartey, Z. Mengesha, C. Touns, J. R. Rickford, D. Jurafsky, and S. Goel, "Racial disparities in automated speech recognition," *Proceedings of the national academy of sciences*, vol. 117, no. 14, pp. 7684–7689, 2020.
- [7] C. Harris, C. Mgbahurike, N. Kumar, and D. Yang, "Modeling gender and dialect bias in automatic speech recognition," in *Findings of the Association for Computational Linguistics: EMNLP 2024*, 2024, pp. 15 166–15 184.
- [8] M. Jahan, P. Mazumdar, T. Thebaud, M. Hasegawa-Johnson, J. Villalba, N. Dehak, and L. Moro-Velazquez, "Unveiling performance bias in asr systems: A study on gender, age, accent, and more," in *ICASSP 2025-2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2025, pp. 1–5.
- [9] S. Joshi, E. Ittan George, T. Javed, K. Bhogale, N. Narasimhan, and M. M. Khapra, "Recognizing Every Voice: Towards Inclusive ASR for Rural Bhojpuri Women," in *Interspeech 2025*, 2025, pp. 4243–4247.
- [10] A. DiChristofano, H. Shuster, S. Chandra, and N. Patwari, "Global performance disparities between english-language accents in automatic speech recognition," *arXiv preprint arXiv:2208.01157*, 2022.
- [11] Z. Liu, I.-E. Veliche, and F. Peng, "Model-based approach for measuring the fairness in asr," in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 6532–6536.
- [12] L. Sari, M. Hasegawa-Johnson, and C. D. Yoo, "Counterfactually fair automatic speech recognition," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 3515–3525, 2021.
- [13] K. L. Lee and S. Watanabe, "The ml-superb 2.0 challenge: Towards inclusive asr benchmarking for all language varieties," 2025.
- [14] A. Radford, J. W. Kim, T. Xu, G. Brockman, C. McLeavey, and I. Sutskever, "Robust speech recognition via large-scale weak supervision," in *ICML*, 2023, pp. 28 492–28 518.
- [15] L. Barrault, Y.-A. Chung, M. C. Meglioli, D. Dale, N. Dong, P.-A. Duquenne, H. Elshar, H. Gong, K. Heffernan, J. Hoffman *et al.*, "Seamlessm4t-massively multilingual & multimodal machine translation," *arXiv preprint arXiv:2308.11596*, 2023.
- [16] M. Donini, L. Oneto, S. Ben-David, J. S. Shawe-Taylor, and M. Pontil, "Empirical risk minimization under fairness constraints," *Advances in neural information processing systems*, vol. 31, 2018.
- [17] M. Pezeshki, O. Kaba, Y. Bengio, A. C. Courville, D. Precup, and G. Lajoie, "Gradient starvation: a learning proclivity in neural networks," *Advances in Neural Information Processing Systems*, vol. 34, pp. 1256–1272, 2021.
- [18] S. Sagawa*, P. W. Koh*, T. B. Hashimoto, and P. Liang, "Distributionally robust neural networks," in *ICLR*, 2020.
- [19] E. Rosenfeld, P. K. Ravikumar, and A. Risteski, "The risks of invariant risk minimization," in *ICLR*, 2021.
- [20] R. Sanabria, N. Bogoychev, N. Markl, A. Carmantini, O. Klejch, and P. Bell, "The edinburgh international accents of english corpus: Towards the democratization of english asr," in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023, pp. 1–5.
- [21] N. Das, S. Bodapati, M. Sunkara, S. Srinivasan, and D. H. Chau, "Best of both worlds: Robust accented speech recognition with adversarial transfer learning," in *Proc. Interspeech 2021*, 2021, pp. 1314–1318.
- [22] I.-E. Veliche, Z. Huang, V. Ayyat Kochaniyan, F. Peng, O. Kalinli, and M. L. Seltzer, "Towards measuring fairness in speech recognition: Fair-speech dataset," in *Proc. Interspeech 2024*, 2024, pp. 1385–1389.