# Bridging the Language Gap: Synthetic Voice Diversity via Latent Mixup for Equitable Speech Recognition

Wesley Bian [1]  Xiaofeng Lin [1]  Guang Cheng [1]

## Abstract

Modern machine learning models for audio tasks often exhibit superior performance on English and other well-resourced languages, primarily due to the abundance of available training data. This disparity leads to an unfair performance gap for low-resource languages, where data collection is both challenging and costly. In this work, we introduce a novel data augmentation technique for speech corpora designed to mitigate this gap. Through comprehensive experiments, we demonstrate that our method significantly improves the performance of automatic speech recognition systems on low-resource languages. Furthermore, we show that our approach outperforms existing augmentation strategies, offering a practical solution for enhancing speech technology in underrepresented linguistic communities.

## 1. Introduction and Related Work

Automatic Speech Recognition (ASR) and other voice-related machine learning technologies have made remarkable advances in recent years, largely propelled by the abundance of data available for English. However, this data imbalance has led to a significant diversity gap: ASR systems perform well for English but often struggle with the world's other 7,000+ languages, perpetuating inequities in access to advanced speech technologies. Addressing this disparity is essential to ensure that the benefits of modern machine learning are distributed equitably across all linguistic communities.

Data augmentation for speech recognition and synthesis has been explored at multiple levels of granularity, including *waveform*, *acoustic feature*, and *latent representation*. Classic perturbative methods—additive noise, speed or

tempo modification, vocal–tract–length warping, and channel convolution—improve robustness, but do not explicitly enhance the diversity of speaker characteristics represented in the dataset. SpecAugment (Park et al., 2019), which introduced time-warping and frequency/time masking on log-Mel spectrograms, remains a strong baseline.

Recent work has shown that interpolation in latent spaces can yield smoother decision boundaries and improved generalization. Mixup (Zhang et al., 2018) and Manifold Mixup (Verma et al., 2019) demonstrate the benefits of linear interpolation between hidden states in vision models, while MixRep (Xie & Hansen, 2023) adapts this approach to ASR by mixing encoder activations, achieving gains on low-resource English corpora. Latent Filling (Bae et al., 2024) applies interpolation to speaker embeddings in zero-shot TTS, enhancing similarity without requiring additional data collection.

Style-based augmentation, particularly through voice conversion models that disentangle speaker characteristics from linguistic content, has further expanded the potential for generating synthetic diversity. Models such as CycleGAN-VC and StarGAN-VC enable non-parallel, many-to-many voice transformations, though audible artifacts in generated audio can limit their downstream utility (Tao et al., 2024). Despite their use in augmentation pipelines, there remains room for improvement.

To the best of our knowledge, no prior work has explored the application of mixup *within* the latent code space of style encoders. We propose LATENTVOICEMIX, a method that operates on this intermediate representation, preserving phonetic structure while expanding the latent convex hull associated with each language. Empirically, we demonstrate that mixup in the style-encoder space yields superior performance compared to existing augmentation methods, controlling for the amount of synthetic audio generated. By bridging latent interpolation theory with codec-level modeling, our approach introduces the first fairness-oriented synthetic data generator at the style layer and provides new evidence that latent convexity is critical for multilingual speech learning.

---
[*]Equal contribution  [1]University of California Los Angeles, Department of Statistics, Los Angeles, United States of America. Correspondence to: Xiaofeng Lin <bernardo1998@g.ucla.edu>, Wesley Bian <wbian@g.ucla.edu>.

## 2. Methodology

We adapt the Diff-HierVC voice conversion model (Choi et al., 2023) to generate synthetic speech data with novel speaker characteristics while preserving the original linguistic content. The core of our approach is to disentangle and manipulate the speaker timbre and linguistic information in audio samples, leveraging a diffusion-based architecture for high-fidelity voice conversion.

### 2.1. Voice Conversion Model Overview

The Diff-HierVC model separates an input audio file into two distinct representations: (1) the linguistic content, corresponding to the words spoken, and (2) the speaker timbre, which encapsulates the unique, non-linguistic characteristics of a person's voice. Speaker timbre refers to the qualities that make a voice unique, independent of linguistic factors such as accent or language. The model enables the recombination of the linguistic encoding from a source audio file with the speaker timbre encoding from a different, target audio file, thereby synthesizing speech that retains the content of the source but adopts the vocal characteristics of the target speaker.

### 2.2. Data Augmentation Procedure

Our augmentation pipeline proceeds as follows:

1. **Audio Cleaning:** All audio files in the input dataset are first denoised using the `noisereduce` Python package (Sainburg et al., 2020) to ensure high-quality input for subsequent processing.

2. **Speaker Timbre Extraction and Storage:** For each audio file in the corpus, we apply the encoding module of Diff-HierVC to extract a fixed-length speaker timbre representation. Specifically, the style encoder produces a 255-dimensional vector that characterizes the unique, time-invariant vocal attributes of each speaker, independent of linguistic content (Choi et al., 2023). These timbre vectors are systematically stored on the filesystem, enabling efficient retrieval and reuse throughout the augmentation pipeline.

3. **Source Selection:** For each data point, the audio file is designated as the *source*, providing the linguistic content for augmentation.

4. **Target Selection:** A separate audio file, spoken by a different speaker, is randomly selected from the dataset to serve as the *target*.

5. **Mixup Timbre Selection:** An additional, pre-saved speaker timbre, distinct from both the source and target speakers, is randomly selected to facilitate mixup.

6. **Voice Conversion with Mixup:** The source linguistic encoding, obtained from the Diff-HierVC model, is combined with a convex combination of the target and mixup speaker timbres. Specifically, let $\mathbf{t}_{\text{target}}$ and $\mathbf{t}_{\text{mixup}}$ denote the timbre vectors of the target and mixup speakers, respectively. The mixed timbre vector is computed as

$$\mathbf{t}_{\text{mixed}} = \lambda\,\mathbf{t}_{\text{target}} + (1 - \lambda)\,\mathbf{t}_{\text{mixup}},$$

where $\lambda \sim \text{Beta}(\alpha, \beta)$ with $\alpha = 0.5$ and $\beta = 0.5$. The model then generates a synthetic audio file that retains the original linguistic content but exhibits a novel, realistic-sounding speaker timbre.

7. **Post-processing:** The synthesized audio is further cleaned using `noisereduce` to remove any residual artifacts.

8. **Transcript Assignment:** The transcript associated with the synthetic audio is inherited directly from the original source file, as the linguistic content remains unchanged.

### 2.3. Benefits and Impact

This augmentation strategy enables the creation of an expanded and more diverse voice corpus without the need for additional data collection. By generating realistic synthetic voices with preserved linguistic content, our method significantly enhances the diversity of training data available for Automatic Speech Recognition (ASR) systems. Empirical results demonstrate that incorporating this augmented data leads to meaningful improvements in ASR model performance, particularly for underrepresented speaker profiles.

## 3. Experiments and Results

We conducted a series of experiments to evaluate the effectiveness of our proposed mixup augmentation technique for improving Automatic Speech Recognition (ASR) performance, and to compare its impact with other established augmentation methods. All experiments report Word Error Rate (WER), a standard metric for ASR evaluation in which lower values indicate better performance.

### 3.1. Datasets

We evaluated the effectiveness of our mixup augmentation technique using three speech corpora: a Wolof dataset (Gauthier et al., 2016), the CSTR VCTK corpus for English (Yamagishi et al., 2019), and the an4 dataset for additional low-resource experimentation (University, 1991). The Wolof corpus consists of 16 hours of transcribed speech from 14 speakers, representing a low-resource language spoken in West Africa. The VCTK corpus provides approximately

44 hours of English speech from 109 speakers with diverse accents, while the an4 dataset is a small English corpus commonly used for benchmarking ASR pipelines in resource-constrained settings.

## 3.2. ASR Models

To rigorously assess the impact of our augmentation method, we trained and evaluated two widely adopted automatic speech recognition (ASR) frameworks: Whisper (Radford et al., 2022) and NVIDIA NeMo (Harper et al., 2019). These models were selected due to their prevalence in both academic and industrial ASR research, as well as their support for multilingual and low-resource scenarios. This experimental design enables us to demonstrate the practical benefits and generalizability of our approach across diverse languages and model architectures.

## 3.3. Alternative Baseline Augmentation Methods

To assess the effectiveness of our proposed mixup augmentation strategy, we conducted comparative experiments against several widely used audio data augmentation techniques. *Waveform augmentation* directly manipulates raw audio signals using time-stretching, amplitude scaling, and pitch shifting, applied randomly to each file with the `audiomentations` library (Jordal et al., 2024). *Spectrogram augmentation*, as in SpecAugment (Park et al., 2019), increases data diversity by masking random frequency bands or time intervals in the spectrogram representation. *Voice conversion augmentation* generates new samples by transferring the linguistic content of an audio file to the vocal characteristics of a different speaker (Zhou et al., 2024), thereby increasing speaker diversity while preserving the original transcript.

## 3.4. AN4: Low-Resource English

We first trained the NVIDIA NeMo ASR model from scratch for 50 epochs on the AN4 dataset. We compared four training regimes: no augmentation, waveform augmentation, voice conversion augmentation, and our mixup augmentation. For all augmentation methods, the dataset size was increased by 33%. Table 1 shows that mixup augmentation significantly improves the performance of this model, and is superior to traditional waveform augmentation.

## 3.5. Wolof vs. English: Addressing Language Bias

To simulate a realistic multilingual training scenario, we constructed a dataset by randomly sampling 8 hours of Wolof speech from the original corpus, ensuring equal representation from each speaker. For English, we similarly sampled 24 hours of data from the VCTK corpus. The NeMo ASR model was then trained for 50 epochs on this com-

*Table 1.* ASR Performance (WER) on AN4 with Different Augmentation Methods

| Training Data | Augmentation | WER |
|---|---|---|
| AN4 only | None | 0.785 |
| AN4 + 33% | Waveform | 0.436 |
| AN4 + 33% | Voice conversion | 0.424 |
| AN4 + 33% | Mixup | **0.339** |

bined dataset. As shown in Table 2, the model achieved substantially better performance on English than on Wolof, highlighting a persistent bias toward the higher-resource language.

We then augmented the Wolof portion of the dataset with an additional 16 hours of synthetic data generated using our mixup augmentation technique, and repeated the training. The results demonstrate a substantial reduction in the performance gap between Wolof and English, indicating improved fairness and inclusivity.

*Table 2.* WER for Multilingual ASR on Wolof and English

| Training Data | Wolof | English | Gap |
|---|---|---|---|
| 8h Wolof + 24h English | 0.796 | 0.562 | 0.234 |
| 24h Wolof (aug.) + 24h English | 0.725 | 0.550 | **0.175** |

## 3.6. Finetuning Whisper on Wolof: Comparison with Other Augmentation Methods

In practice, many automatic speech recognition (ASR) systems are developed by fine-tuning large pretrained models, such as Whisper by OpenAI, rather than training from scratch. To assess the impact of our augmentation methods in this setting, we fine-tuned the `whisper-tiny` model on 8 hours of original Wolof data for 4 epochs. For each original sample, two synthetic samples were generated using various augmentation techniques, resulting in a tripled dataset size. We also evaluated the perceptual quality of the augmented data using the average SpeechMOS score (Huang et al., 2022), or, in the case of no augmentation, the original data. SpeechMOS is a neural model that predicts human-perceived speech quality on a scale from 1 to 5. For spectrogram augmentation, SpeechMOS was not computed, as this method does not generate waveform outputs.

Table 3 summarizes the WER achieved with different augmentation strategies. Our mixup augmentation method consistently outperformed spectrogram augmentation, waveform augmentation, and conventional voice conversion augmentation, yielding the lowest WER.

*Table 3.* WER for Whisper Finetuned on Wolof with Different Augmentation Methods

| Augmentation Method | WER | SpeechMOS |
|---|---|---|
| None | 0.283 | 2.661 |
| Spectrogram Augmentation | 0.242 | n/a |
| Waveform Augmentation | 0.217 | 2.117 |
| Voice Conversion Augmentation | 0.215 | 2.710 |
| Mixup Augmentation (proposed) | **0.202** | 2.243 |

## 4. Analysis of Method

### 4.1. Ablation Study on Mixup Augmentation

To understand the contribution of individual components within our mixup augmentation algorithm, we conducted an ablation study using 8 hours of original Wolof data for fine-tuning Whisper. The effectiveness of each variant was assessed using word error rate (WER) on the ASR task. Table 4 summarizes the results for various configurations.

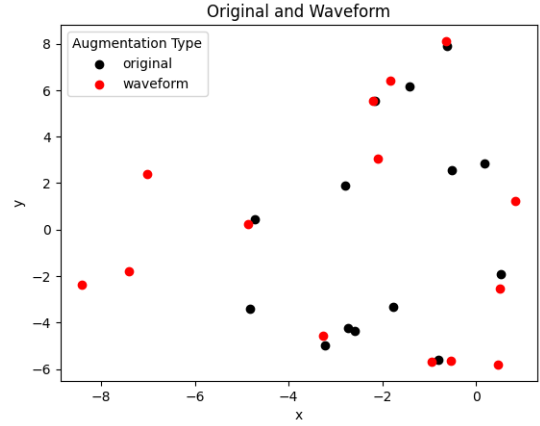*Table 4.* Ablation study of mixup augmentation variants on Wolof ASR (WER; lower is better).

| Setting | WER |
|---|---|
| No Post-denoising, Source=Target (8h) | 0.235 |
| No Post-denoising, Source=Target (16h) | 0.221 |
| Mixup w/ 3 Speaker Timbres (16h) | 0.221 |
| No Post-denoising (16h) | 0.214 |
| Proposed Mixup (16h) | **0.202** |

The results indicate that the full mixup algorithm, including post-denoising, achieves the lowest WER. Omitting the post-denoising step or altering the mixup configuration—such as using three timbres or setting the source and target to the same audio file—consistently results in higher error rates. These findings underscore the importance of post-denoising and careful design of the mixup process for optimal augmentation performance.
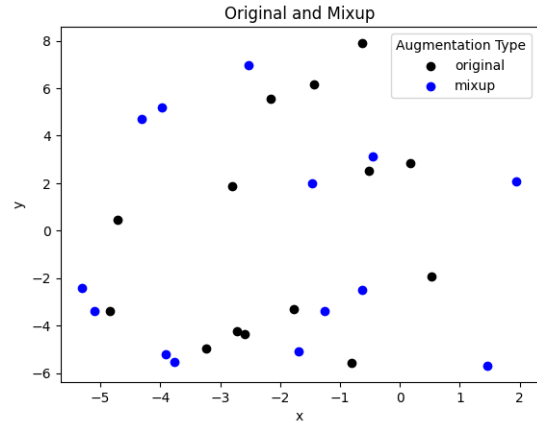
### 4.2. Analysis of Speaker Timbre Distributions

To further investigate the qualitative differences between augmentation methods, we performed principal component analysis (PCA) on the speaker timbre vectors of the 14 speakers in the Wolof corpus. We compared these to timbre vectors extracted from synthetic samples generated by both mixup and traditional waveform augmentation. As illustrated in Figure 4.2, the speaker timbres produced by mixup are distributed more closely to those of the original speakers, whereas timbres from waveform augmentation exhibit greater variance and tend to lie outside the distribution of real speaker timbres. These results suggest that mixup augmentation generates synthetic data that not only increases training diversity but also more faithfully preserves the underlying structure of the original speaker timbre distribution. This property may explain the superior performance of mixup augmentation compared to traditional waveform augmentation.





## 5. Conclusion

We introduced a novel data augmentation technique that applies mixup in the latent space of voice conversion models, and demonstrated its effectiveness in enhancing Automatic Speech Recognition (ASR) performance. Our experiments show that this method outperforms traditional augmentation techniques, particularly for low-resource languages. These findings indicate that significant gains in ASR for underrepresented languages can be achieved without extensive data collection, promoting broader access to advanced speech technologies.

## Impact Statement

This work advances machine learning for audio in low-resource languages by reducing bias toward well-resourced languages. Our approach promotes equitable access to speech technology, enabling speakers of all languages to benefit from progress in machine learning systems.

# References

Bae, J.-S., Lee, J. Y., Lee, J.-H., Mun, S., Kang, T., Cho, H.-Y., and Kim, C. Latent filling: Latent space data augmentation for zero-shot speech synthesis. In *ICASSP 2024 Workshop on Speech and Language Processing*, 2024.

Choi, H.-Y., Lee, J.-H., Choi, H.-S., Byun, K., Lee, K., Hwang, M.-J., Kim, M., and Lee, K. Diff-hiervc: Diffusion-based hierarchical voice conversion with robust pitch generation and masked prior for zero-shot speaker adaptation. *arXiv preprint arXiv:2311.04693*, 2023.

Gauthier, E., Besacier, L., Voisin, S., Melese, M., and Elingui, U. P. Collecting resources in sub-saharan african languages for automatic speech recognition: a case study of wolof. *LREC*, 2016.

Harper, E., Majumdar, S., Kuchaiev, O., Li, J., Zhang, Y., Bakhturina, E., Noroozi, V., Subramanian, S., Koluguri, N., Huang, J., Ginsburg, B., Gadde, R., Nguyen, H., Leary, R., Cohen, J., Tumanov, D., Agarwal, A., Rosenzweig, C., Tovstogan, P., Lavrukhin, V., and Wang, Y. Nemo: a toolkit for conversational ai and large language models, 2019. URL https://nvidia.github.io/NeMo/. Version 1.0.0. https://github.com/NVIDIA/NeMo.

Huang, G., Wang, Q., Zhang, Y., et al. Speechmos: A universal non-intrusive mos predictor for synthetic speech. *arXiv preprint arXiv:2204.02152*, 2022.

Jordal, I., Tamazian, A., Dhyani, T., Chourdakis, E. T., Karpov, N., Landschoot, C., Angonin, C., Sarioglu, O., BakerBunker, kvilouras, Çoban, E. B., Gritskevich, E., Mirus, F., Lee, J.-Y., Choi, K., Killingberg, L., Marvin-Lvn, SolomidHero, Alumäe, T., and Solovyev, R. audiomentations: v0.38.0, December 2024. URL https://github.com/iver56/audiomentations.

Park, D. S., Chan, W., Zhang, Y., Chiu, C.-C., Zoph, B., Cubuk, E. D., and Le, Q. V. Specaugment: A simple data augmentation method for automatic speech recognition. In *Proceedings of Interspeech*, pp. 2613–2617, 2019.

Radford, A., Kim, J. W., Xu, T., Brockman, G., McLeavey, C., and Sutskever, I. Robust speech recognition via large-scale weak supervision. *arXiv preprint arXiv:2212.04356*, 2022. URL https://arxiv.org/abs/2212.04356.

Sainburg, T., Thielk, M., and Gentner, T. Q. Finding, visualizing, and quantifying latent structure across diverse animal vocal repertoires. *PLoS computational biology*, 16(10):e1008228, 2020.

Tao, R., Zhang, Y., Gong, Y., and Ling, Z. Voice conversion augmentation for speaker recognition on defective datasets. *arXiv preprint arXiv:2404.00863*, 2024.

University, C. M. An4 speech dataset. https://github.com/cmusphinx/an4, 1991. Recorded at Carnegie Mellon University.

Verma, V., Lamb, A., Beckham, C., Najafi, A., Mitliagkas, I., Bengio, Y., and Courville, A. Manifold mixup: Better representations by interpolating hidden states. In *Proceedings of the 36th International Conference on Machine Learning*, pp. 6438–6448, 2019.

Xie, J. and Hansen, J. H. L. Mixrep: Hidden representation mixup for low-resource speech recognition. In *Proceedings of Interspeech*, pp. 1304–1308, 2023. doi: 10.21437/Interspeech.2023-1216.

Yamagishi, J., Veaux, C., and MacDonald, K. 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. URL https://datashare.ed.ac.uk/handle/10283/3443.

Zhang, H., Cisse, M., Dauphin, Y. N., and Lopez-Paz, D. mixup: Beyond empirical risk minimization. In *International Conference on Learning Representations*, 2018.

Zhou, Z., Xu, S., Yin, S., Li, L., and Wang, D. A comprehensive investigation on speaker augmentation for speaker recognition. *arXiv preprint arXiv:2406.07421*, 2024.