

# ENHANCING NON-NATIVE ACCENT RECOGNITION THROUGH A COMBINATION OF SPEAKER EMBEDDINGS, PROSODIC AND VOCAL SPEECH FEATURES

Anup Bera and Aanchal Agarwal

Accenture India

## **ABSTRACT**

*The transcription accuracy of automatic speech recognition (ASR) system may suffer when recognizing accented speech. The resulting bias in ASR system towards a specific accent due to under representation of that accent in the training dataset. Accent recognition of existing speech samples can help with the preparation of the training datasets, which is an important step toward closing the accent gap and eliminating biases in ASR system. For that we built a system to recognize accent from spoken speech data. In this study, we have explored some prosodic and vocal speech features as well as speaker embeddings for accent recognition on our custom English speech data that covers speakers from around the world with varying accents. We demonstrate that our selected speech features are more effective in recognizing non-native accents. Additionally, we experimented with a hierarchical classification model for multi-level accent classification. To establish an accent hierarchy, we employed a bottom-up approach, combining regional accents and categorizing them as either native or non-native at the top level. Furthermore, we conducted a comparative study between flat classification and hierarchical classification using the accent hierarchy structure.*

## **KEYWORDS**

*Automatic speech recognition; accent recognition; Hierarchical classification; MLP Classifier; Speech Features; speaker embedding; native and non-native accents*

## **1. INTRODUCTION**

The advent of automatic speech recognition (ASR) system and its wide adaptation in industry across different applications such as voice bot, voice search, mobile phones, home appliances etc., demands high transcription accuracy. The performance of some of the ASR systems is not up to the desirable mark for accented speech, especially for non-native speakers. The poor performance for a particular accented speech may be due to poor representation of that accented speech in the training dataset. So, identifying accent gap and enriching training dataset by including different accented speech data helps to improve ASR model's overall performance. To address this problem, it is necessary to recognize spoken accent of the speech data and using that information create a well-balanced speech training dataset. We built a supervised classification model for accent recognition. In this paper we examine the comparative study of the flat classification and hierarchical classification for accent recognition for English speech data. The classification models are trained using input features such as prosodic and vocal speech features as well as speaker embeddings of the speech data. Here we also demonstrate that overall model accuracy will be improved if prosodic and vocal speech features are being used along with speaker embedding vectors. It has also been observed that, our speech features enhance model

performance towards non-native accent recognition. As our accent recognition experiment is not limited for few selected accents or specific to any geolocation, we have prepared a custom dataset which covers speakers from around the world with varying accents.

### **1.1. Related Work**

Most of the previous works [1], used MFCC as the input features and chose neural network model for the accent prediction. A comparative study on various speech features, such as, MFCC, Spectrogram, Chromagram, Spectral Centroid, and Spectral Roll-off has been conducted by [2] using 2 layers CNN model. And they included that MFCC has given best accuracy on 5 accents, namely, Arabic, English, French, Mandarin and Spanish, with 48.4% accuracy. Also, Arlo Faria [3] explored accent classification for native and non-native speaker using acoustic and lexical features, achieving accuracy 84.5%. The L. M. K. Sheng [4] explored deep learning approach for accent classification on 3 accents, namely, English, Chinese, Korea, and achieved 87.6% accuracy. There are some other works which jointly studied the relationship between ASR with accent identification and use that information to select the accented speech [5] or improve accuracy of both [6] or use accent information to improve ASR accuracy [7]. So far, no experiment done using speaker embeddings as input features for accent recognition. As speaker embeddings carry some salient information, such as speaker speaking style, age, dialect, accent etc., it intrigued us to experiment with this feature. Further we appended some prosodic and vocal features to embedding vectors to understand the accent recognition model performance on the combined features.

## **2. RESEARCH METHODOLOGY**

### **2.1. Dataset Preparation**

We have prepared a custom dataset by combining multiple open datasets that provides sufficient representation of various accents around the world.

#### **2.1.1. Dataset Description**

Our experiment is to recognize accent for an English speech which is spoken anywhere in the world. In order to achieve that, we have prepared a dataset which are combination of various speech data spoken across the geographic regions. To conduct our experiments, we used 5 open-source speech datasets: Kaggle Common Voice [8], George Mason University Speech Accent Archive [9], CSTR VCTK Corpus [10], NISP [11] and TIMIT [12] datasets. We have combined these 5 datasets to train our model. For this experiment we have used a total of 22.5 hours of speech data by under-sampling some accent groups to make dataset more balanced and comprehensive. The detail of the dataset preparation approach is given in the paper [13].

#### **2.1.2. Data Processing**

Each dataset has different file format, sampling rate, encoding, and file duration. For example, TIMIT dataset has files with duration mostly less than 5 secs and having 16 KHz sampling rate with wav format. Whereas Common Voice data have 48 KHz sampling rate and in mp3 format. The VCTK data are in flac format with 48 kHz sampling rate. Also, in some dataset there are some language data other than English. As this study focus on English speech data, it is required to filter out only for English data from the original dataset.

### 2.1.3. Data Standardization

As the above 5 datasets have different file format, sampling rate, bit rate, encoding, and file duration, to create a single training dataset out of these, it is necessary to standardize all speech data. Regardless of the original file format, sampling rate and bit rate of each dataset, a pipeline was implemented to standardize all the data. Also, speech file duration needs to be considered during data preparation. After some initial experiments, we found that a file with min duration of 5 seconds and max duration of 20 seconds give optimum result. So, we filtered out only those files having duration more than 5 seconds and clipped all the files which are having duration more than 20 seconds. All the data are standardized to following format:

- File format: Wav
- Sampling rate: 16KHz
- Bits per sample: 16
- Encoding: PCM\_S
- Channel: single
- File max duration: 20secs
- File min duration: 5secs

## 2.2. Accent Label preparation

Wherever the accent labels are available, we have used those as it is. In case the accent labels directly are not available, we have used indirect information, such as speaker first language, mother tongue, speaker coming from which region etc, to assign accent label. For example, if metadata of the dataset mentions that speaker's first language is Hindi, then the accent of the data is assigned as Indian. Similarly, if it is mentioned that speaker's first language as French, then accent of the speech data assigned as French.

## 2.3. Hierarchical Accent Preparation

Now we have accent information for each speech data in the combined dataset. As some of the accents have shared characteristics, those accents are grouped together and created a hierarchical accent tree structure using bottom-up approach. We have created one variable, called, "geolocation\_accent" by grouping all the accents which are originated from similar geographic location. For example, Arabic, Farsi, Hebrew, Kurdish accents are spoken in middle east geographic location. All these accents from middle east are grouped together and map to "Middle East" label in "geolocation\_accent" variable. The purpose to create this type of geolocation accent purely basis on needs of business use case. Mostly companies are working in a specific geographic location, and they have operational unit to serve that geographical area. In that scenario it is required to recognize the speaker geolocation\_accent and based on that companies can take necessary action. The total number of geolocation accent labels is 9; these are 'Africa', 'Australasia', 'British', 'East Asia', 'Europe', 'Middle East', 'North America', 'South Asia', 'South East Asia'. For detail of this grouping and mapping of accents, please refer the Table 11 in Appendix.

Further, for all those geolocation\_accent labels in which first language is English, are combined to "Native" accent category and rests are combined to "Non-Native" accent category where first language is other than English. A new variable is created, namely, "binary\_accent" comprising with "native" and "non-native" accent labels. Here "North American", "Australasia" and "British" are grouped to "Native" label of "binary\_accent" and "Africa", "East Asia", "Europe", "Middle East", "South Asia", "South East Asia" labels are grouped to "Non-Native" label of "binary\_accent". So, we have prepared 2 levels hierarchical accent structure. At the top, it is

“binary\_accent” with “Native” and “Non-Native” labels, and “geolocation\_accent” comprising 9 labels is in 2<sup>nd</sup> level.

The accent speech file distributions and counts of each label of the “binary\_accent” and “geolocation\_accent” are shown in Table 1 and Table 2.

Table 1: native and non-native accent distribution in binary\_accent variable

<b>binary_accent</b>	<b>File counts</b>	<b>Distribution percentage</b>	<b>Total duration in seconds</b>
Native	2998	44%	69259
Non-Native	3802	56%	11395

Table 2: The distribution of the different geolocation\_accent categories.

<b>geolocation_accent</b>	<b>Filecounts</b>	<b>File distribution percentage</b>	<b>Total duration in seconds</b>
Africa	733	10.78	5788.7
Australasia	998	14.68	3838.1
British	1000	14.70	3716.8
East Asia	406	5.97	3840.4
Europe	978	14.38	8211.2
Middle East	251	3.69	27197.1
North America	1000	14.70	7394.6
South Asia	1000	14.70	12003.3
South East Asia	434	6.38	8663.7

## 2.4. Features Preparation

For accent recognition, we use several prosodic and vocal speech features as well as speaker embeddings as model inputs. Hereafter we will discuss the methods to create various speech features.

### 2.4.1. Speaker Embeddings

We have chosen freely available speaker embedding pre-trained model which is built by Real-Time-Voice-Cloning [14] team. The model generates embedding vector having size of 512 for a given speaker speech. The embedding vectors will be remained almost same for all speech data from same speaker. The team has implemented the model [15] using the paper [16] and trained the model using the 3 datasets, Libri speech [17], Voxceleb1 [18], Voxceleb2 [19]. The team has used this model for their speaker recognition project.

### 2.4.2. Speech Features

Initially, we have started with more than 20 various prosodic and vocal speech features. But after some analysis, we have identified 7 prosodic and vocal speech features, namely, harmonicity to noise, pitch mean, speech rate, number of pauses, speaking ratio, energy, and pitch frequency mean, have shown high correlation, either positive or negative, with “geolocation\_accent” accent

labels. The heatmap Figure 1 is showing the correlation values of each of the features with the “geolocation\_accent” labels. The correlation has been calculated using Pearson approach. To extract those speech features, we have used librosa [20], parselmouth [21] and prosody [22] libraries. The Table 3 provides details of each of the 7 speech features, such as, feature name, unit of measurement, description, computed values, and library used.

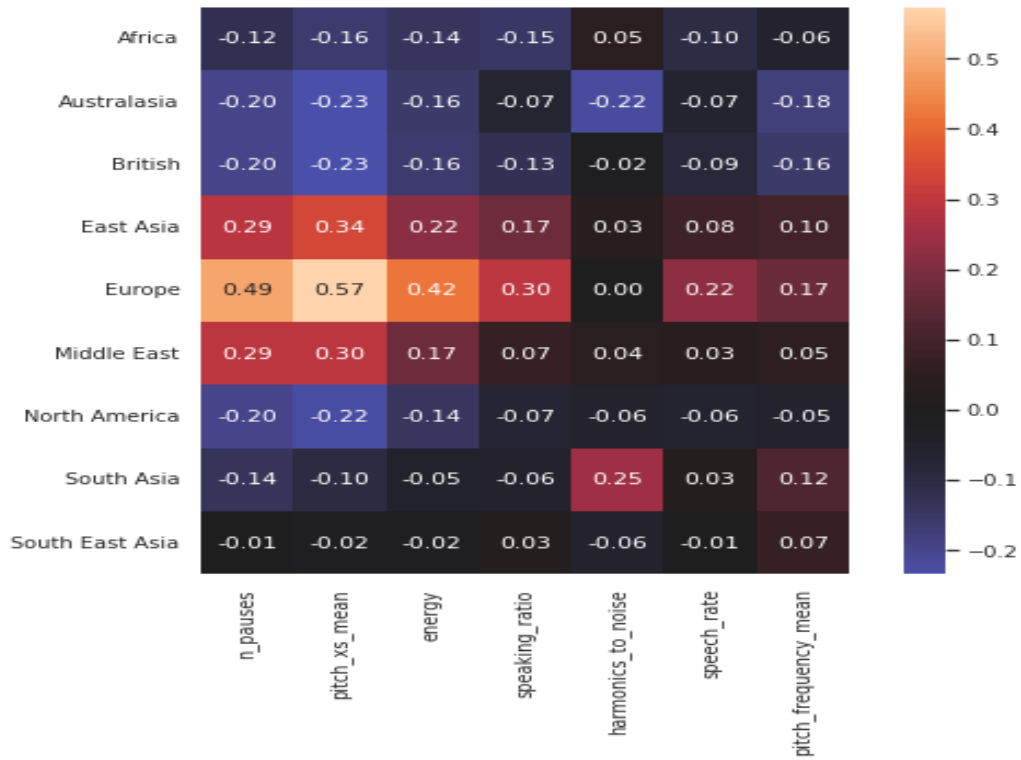


Figure 1: Correlation heatmap among speech features and geolocation accents

Table 3: Speech features extraction methods

Feature Name	Unit of Measurement	Description	Computed Values	Library and Function Used
harmonics_to_noise	Decibel (dB)	Ratio between f0 and noise components, which indirectly correlates with perceived aspiration. This may be due to reducing laryngeal muscle tension resulting in a more open, turbulent glottis.	Harmonicity Mean	Library: parselmouth Function: sound.to_harmonicity()
pitch_xs_mean	dB	Pitch is a measure of how high or low something sounds and is related to the speed of the vibrations that produce the sound	Pitch XS – Mean	Library: parselmouth Function: sound.to_pitch()
speech_rate	Utterances/Seconds	Number of speech utterances per second over the duration of the speech sample (including pauses)	Speech Rate	Library: Myprosody Function: myspsr()
n_pauses		Number of Pauses taken in a piece of audio	Number of Pauses	Library: MyProsody Function: mysppaus()
speaking_ratio		The measure of ratio between speaking duration and total speaking duration	Speaking Ratio	Library: MyProsody Function: myspbala()
energy	db	Measured as the mean-squared central difference across frames and may correlate with motor coordination.	Energy	Library: parselmouth Function: sound.get_energy()
pitch_frequency_mean	db	Pitch is a measure of how high or low something sounds and is related to the speed of the vibrations that produce the sound. Here Pitch Frequency gives the pitch of the fragment	Pitch Frequency Mean	Library: parselmouth Function: pitch =sound.to_pitch() pitch.selected_array['frequency'].mean()

## 2.5. Experimental Design

Experiment 1: Binary classification on the target variable “binary\_accent”

Exp 1.1 Binary classification between “Native” and “Non-Native” labels of the target variable “binary\_accent” using speaker embedding vectors as input features

Exp 1.2 Binary classification between “Native” and “Non-Native” labels of the target variable “binary\_accent” using input feature as combination of speech features and speaker embedding vectors.

Experiment 2: Multi-label classification on the target variable “geolocation\_accent”

Exp 2.1 Multi-label classification for 9 “geolocation\_accent” labels using input feature as speaker embedding vectors.

Exp 2.2 Multi-label classification of 9 “geolocation\_accent” labels using input features as combination of speech features and speaker embedding vectors.

Experiment 3: Hierarchical classification of the accents using HiClass[23].

Exp 3.1 Hierarchical classification on hierarchical structure of “binary\_accent” and “geolocation\_accent” using input feature as speaker embedding vectors.

Exp 3.2 Hierarchical classification on hierarchical structure of “binary\_accent” and “geolocation\_accent” using input feature as combination of speaker embedding vectors and speech features.

## 2.6. Train and Test Dataset Preparation

For all the experiments, the data has been split in 3:1 ratio for train and test dataset with stratified on target variable “binary\_accent” so that Native and Non-Native accents are remain in equal ratio in both the train and test dataset. Further the dataset also stratified on “geolocation\_accent” classes. All the geolocation accents are distributed in equal ratio between train and test dataset. The accent file counts distribution in the dataset is shown in the Table 4.

Table 4: Accent file count distribution in train and test set

dataset/ geolocation_ accents	Australasia	British	North America	Africa	East Asia	Europe	Middle East	South Asia	South East Asia
dataset/ binary_ accent	Native			Non-Native					
train	748	750	750	550	305	733	188	750	325
	2248			2852					
test	250	250	250	183	101	245	63	250	109
	750			950					

## 3. EXPERIMENTS AND RESULTS

### 3.1. Experiment 1

The experiment is carried out on “binary\_accent” as target variable whose classes are “Native” and “Non-Native”.

**Model:** In the first experiment, we used multi-layer perceptron model for binary classification. For this experiment purpose we have used scikit-learn MLP Classifier with the hyperparameters as 1 hidden layer with 100 nodes, relu as activation function and adam as estimator with log-loss as loss function.

**Features:** In this binary classification experiment, first we have used input features as speaker embedding vectors as mentioned in Features Preparation section and for the second experiment, the input features are prepared by appending 7 speech features with the speaker embedding vectors.

**Result:** Through the experiment we have found the optimum model hyperparameters which are mentioned in the Table 5. The accuracy is calculated on the test dataset which is prepared as described in the section 2.6. The results of the 2 separate experiments conclude that after including 7prosodic and vocal speech features along with embedding vectors giving better accuracy. The result is not calculated basis on overall prediction accuracy, rather consider individual class accuracy and then average out on both Native and Non-Native accuracies. It should be noticed from individual accuracy result, which is shown in the Table 6, that Non-Native accent prediction accuracy is improved when 7 speech features are added in the model input feature.

Table 5: Overall accuracy when classification model trained on embedding vectors and combined with speech features

Experiment	Target Variable	Classes	Features	Classification model	Model hyperparameters	Accuracy	F1 score
Exp1.1	binary_accent	Native, Non-Native	Speaker Embedding	MLPClassifier	random_state=1, max_iter=500, shuffle=True	90.71%	91
Exp1.2	binary_accent	Native, Non-Native	Speaker Embedding + speech features	MLPClassifier	random_state=1, max_iter=500, shuffle=True	92.5%	93

Table 6: Individual native and non-native accent accuracy in percentage

binary_accent	Accuracy of Exp 1.1	Accuracy of Exp 1.2
Native	90.3%	90.8%
Non-Native	91.1%	94.2%

### 3.2. Experiment 2

This experiment is carried out on “geolocation\_accent” as target variable which have 9 accent classes.

**Model:** Here the model remains same as for binary classification with only difference in the output layer which has 9 nodes.

**Features:** First experiment is carried out using embedding vectors as input features and then in the second experiment, the speech features are appended to the embedding vectors to create a combined input feature.

**Results:** The model hyperparameters are remain same as in the experiment 1. Here again the final accuracy is calculated based on each of the classes’ accuracy in test dataset and then average out of all individual accuracy. The individual accuracy of each class is shown in the Table 8 and the overall accuracy is given in the Table 7.



Table 7: Overall accuracy when classification model trained on embedding vectors and combined with speech features

Experiment	Target Variable	Classes	Features	Classification model	Model hyperparameters	Accuracy	F1 score
Exp 2.1	geolocation_accent	9 classes	Speaker Embedding	MLPClassifier	number of layers=1, hidden node=100, random_state=1, max_iter=500, shuffle=True	60.24%	60
Exp 2.2	geolocation_accent	9 classes	Speaker Embedding + speech features	MLPClassifier	number of layers=1, hidden node=100, random_state=1, hidden nodes=100, max_iter=500, shuffle=True	62.14%	62.78

Table 8: Individual accuracy for each accent in geolocation\_accent

geolocation_accent	Accuracy of Exp 2.1	Accuracy of Exp 2.2
Africa	63%	62%
Australasia	81%	82%
British	67%	74%
East Asia	34%	40%
Europe	61%	76%
Middle East	25%	17%
North America	70%	75%
South Asia	84%	86%
South East Asia	57%	47%

### 3.3. Experiment 3

Next, we have experimented hierarchical classification on both “binary\_accent” and “geolocation\_accent” variables.

**Model:** For the hierarchical classification, we have used HiClass [23] library. The HiClass is compatible with scikit-learn [24] and support various hierarchical classifiers, such as, Local Classifier per Parent Node, Local Classifier per Node, Local Classifier per Level. In this experiment, the Local Classifier per Parent Node approach is chosen because it consists of training a multi-class classifier for each parent node existing in the hierarchy. This is exactly the current problem statement. First classify a speech data either to Native or Non-Native category and then further classify to the Native geolocation accents or Non-Native geolocation accents. The HiClass needs a base classification model which should be compatible with scikit-learn. To keep all the experiments in comparable, we chose MLP Classifier as base classification model. Features: The features are remained same as with the previous experiments. First experiment is done with speaker embedding vectors as input features and then in next experiment, 7 speech features are combined with the speaker embedding vectors.

**Results:** The base model, MLP Classifier’s hyperparameters remain same as in the previous experiments. In case of hierarchical classification model, the default hyperparameters are used. Here also it can be noticed that after appending speech features, model giving overall better accuracy. In this experiment also, the Non-Native accent as well as non-native geolocation accent recognition accuracies have improved after appending 7 speech features to embedding vectors as input features. Only exception is for “Middle East” accent, where accuracy is decreased. The individual accent accuracy is shown in the Table 10 and overall accuracy is shown in the Table 9 for Native, Non-Native and geolocation accents.

Table 9: Overall accuracy when hierarchical classification model trained on embedding vectors and combined with speech features

Experiments	Number of hierarchies	Classes	Features	Classification base model	Base Model hyperparameters	Hierarchical Classification Model	Accuracy at level 1	F1 score at level 1	Accuracy at Level 2	F1 score at level 2
Exp3.1	2	Level 1: Native, Non-Native Level 2: 9 classes	Speaker Embedding	MLPClassifier	number of layers=1, hidden node=100, random_state=1, max_iter=500, shuffle=True	LocalClassifierPerLevelParentNode	90.7%	91	62.7%	62.7
Exp3.2	2	Level 1: Native, Non-Native Level 2: 9 classes	Speaker Embedding + speech features	MLPClassifier	number of layers=1, hidden node=100, random_state=1, max_iter=500, shuffle=True	LocalClassifierPerLevelParentNode	92.5%	93	63.2%	63.3

Table 10: Individual accuracy from hierarchical classification model for each accent in binary\_accent and geolocation\_accent

binary_accent	Accuracy of Exp 3.1	Accuracy of Exp 3.2
Native	90.3%	90.8%
Non-Native	91.1%	94.2%
geolocation_accent		
Africa	65%	67%
Australasia	87%	84%
British	76%	67%
East Asia	42%	46%
Europe	53%	72%
Middle East	35%	24%
North America	74%	72%
South Asia	82%	85%
South East Asia	51%	52%

#### **4. DISCUSSION**

In case of “Native” and “Non-Native” binary classification, both flat classification from experiment 1 and hierarchical classification from experiment-3 give same result. But in case of “geolocation\_accent” classification, the hierarchical classification in experiment-3 giving better result than flat classification in experiment-2. The above result has shown that the accuracy for “East Asia”, “Middle East” and “South East Asia” are lower than other classes. That we believe due to low representation of those accented speech data in training dataset. Another interesting finding is that, after including selected speech features, model performance has improved towards non-native geolocation accents with exception of “Middle East”.

#### **5. CONCLUSION AND FUTURE WORK**

In this work we have demonstrated four key points. First, using embedding vectors we can get decent accuracy for accent recognition for wide variety of accents. Second, by combining prosodic and vocal speech features along with speaker embedding vectors, the classification result showing improvement. Third, instead of use 2 separate flat classification models for 2 level classification, it is always better to use a single hierarchical classification model. Fourth, our selected 7 prosodic and vocal speech features improve the non-native accent recognition.

The size of training dataset for each geolocation accent is very low. We belief due to lack of sufficient data, the overall accuracy is not so promising. In future, we will include more training data as well as better hierarchical and flat classification model to improve the overall accuracy. Also, our next focus area will be to increase the number of hierarchical levels further down to regional dialects.

#### **Appendix**

The accent mapping table for binary\_accent and geolocation\_accent variables from the original data are shown in Table 11.

Table 11: Accent mapping with geolocation accents and native, non-native accents

Binary Accent	Accent Mapping	Accent
Native	North America	Canadian American New England North Midland South Midland
Native	British	Irish Welsh British Scottish
Native	Australasia	Australian New Zealand
Non-Native	South Asia	Indian Nepali Sinhala Burmese
Non-Native	Africa	Hausa African Amharic Kiswahili South African
Non-Native	Europe	French Spanish Russian German Portuguese
Non-Native	South-East Asia	Malaysia Singapore Philippines Vietnamese
Non-Native	East Asia	Korean Japanese Mandarin Hongkong Cantonese
Non-Native	Middle East	Farsi Arabic Kurdish Hebrew

## REFERENCES

- [1] C. Shih, "Speech Accent Classification," Project Report, Stanford University, Stanford, CA, 2017.
- [2] Y. e. e. Singh, "Features of Speech Audio for Accent Recognition," 2020 International Conference on Artificial Intelligence, IEEE, 2020.
- [3] A. Faria, "Accent Classification for Speech Recognition," International Computer Science Institute, Berkeley CA 94704, USA, p. 285–293, 2005.
- [4] L. M. A. Sheng, "Deep Learning Approach to Accent Classification," roject Report, Stanford University, Stanford, CA, 2017.
- [5] M. & R. M. Najafian, "Automatic accent identification as an analytical tool for accent robust automatic speech recognition.," Speech Communication, 2020.
- [6] K. C. S. & M. L. Deng, "Improving accent identification and accented speech recognition under a framework of self-supervised learning," arXiv preprint arXiv:2109.07349, 2021.

- [7] K. Sreenivasa Rao, "Transfer Accent Identification Learning for Enhancing Speech Emotion Recognition".
- [8] "Kaggle CV," [Online]. Available: <https://www.kaggle.com/datasets/mozillaorg/common-voice>.
- [9] "GMU accent," [Online]. Available: [https://accent.gmu.edu/browse\\_language.php?function=find&language=english](https://accent.gmu.edu/browse_language.php?function=find&language=english).
- [10] J. Yamagishi, "VCTK," [Online]. Available: <https://datashare.ed.ac.uk/handle/10283/3443>.
- [11] S. B. V. D. G. S. K. P. Kalluri, "NISP: A Multi-lingual Multi-accent Dataset for Speaker Profiling," In ICASSP 2021-2021 IEEE International Conference on Acoustics, pp. 6953–6957, IEEE, 2021.
- [12] J. S. Garofolo, "TIMIT," 1993. [Online]. Available: <https://catalog.ldc.upenn.edu/LDC93s1>.
- [13] A. Bera and V. V. Kandasamy, "IMPROVING ROBUSTNESS OF AGE AND GENDER PREDICTION BASED ON CUSTOM SPEECH DATA," in AIRCC, 2022.
- [14] "voice cloning," 2021. [Online]. Available: <https://github.com/CoirentinJ/Real-Time-Voice-Cloning>.
- [15] "speaker encoder," 2021. [Online]. Available: <https://github.com/CoirentinJ/Real-Time-Voice-Cloning/tree/master/encoder>.
- [16] W. Q. R. A. M. I. L. Wan L., "GENERALIZED END-TO-END LOSS FOR SPEAKER VERIFICATION," Retrieved from <https://arxiv.org/pdf/1710.10467.pdf>, 2020.
- [17] C. G. P. D. K. S. Panayotov V., "LIBRISPEECH: AN ASR CORPUS BASED ON PUBLIC DOMAIN AUDIO BOOKS," ICASSP, 2015.
- [18] C. J. S. X. W. Z. A. Nagrani A., "Voxceleb: Large-scale speaker verification in the wild," Computer Science and Language,, 2019.
- [19] N. S. Z. A. Chung J. S., "VoxCeleb2: Deep Speaker Recognition," INTERSPEECH, 2018.
- [20] "librosa," [Online]. Available: <https://librosa.org>.
- [21] "praat," [Online]. Available: <https://parselmouth.readthedocs.io/>.
- [22] S. Shahab, "prosody," [Online]. Available: <https://shahabks.github.io/myprosody/>.
- [23] K. N. R. B. Y. Miranda F. M., "HiClass: a Python library for local hierarchical classification compatible with scikit-learn," Retrieved from <https://arxiv.org/pdf/2112.06560.pdf>, 2022.
- [24] "scikit-learn," [Online]. Available: <https://scikit-learn.org/stable/>.
- [25] A. T. a. S. U. K. Chakraborty, "Voice Recognition Using MFCC Algorithm," International Journal of Innovative Research in Advanced Engineering, vol. 1, no. 10, p. 4, 2014.

#### AUTHORS

**Anup Bera** has done M.Tech from IIT Kharagpur and he has total 17 years of industry experience and currently working in Accenture India



**Aanchal Agarwal** has done Master in Computer Science in 2020 and currently she is working as Machine Learning expert in Accenture India

