# A Survey on Telugu Accent Classification and Conversion

## Mrs. B. Sreelatha[1], Madastham Pranathi[2], K Manisha Reddy[3],J Shirisha[4]

*Professor and Head, Department of CSE (Artificial Intelligence & Machine Learning)[1]*
*IV B. Tech Students, Department of CSE (Artificial Intelligence & Machine Learning)[2,3,4]*
*ACE Engineering College, Hyderabad, Telangana, India*

--------------------------------------------------------------------------------------------------------------------------
--------------------------------------------------------------------------------------------------------------------------

**ABSTRACT**: The diversity of accents within the Telugu language presents a significant challenge in effective communication. This research explores a novel approach to address this issue by proposing a comprehensive system that combines accent classification and conversion using artificial intelligence (AI). The primary objective is to enhance communication between individuals with distinct Telugu accents. The proposed system leverages machine learning techniques, Convolutional Neural Networks (CNNs), to achieve precise Telugu accent classification. By training the model on a diverse dataset encompassing various Telugu accents, the system becomes adept at recognizing and categorizing the subtle nuances in pronunciation, intonation, and speech patterns that characterize different accents. Furthermore, the research introduces the concept of Telugu accent conversion, aiming to transform one Telugu accent into another to facilitate improved understanding among speakers with differing linguistic backgrounds. This involves the integration of speech recognition and text-to-speech modules, allowing the system to transcribe spoken words into text and convert them into a different accent while maintaining naturalness and clarity. The developed system holds promise in bridging communication gaps caused by accent variations, ultimately contributing to more inclusive and effective verbal communication in Telugu.

**Keywords:** Telugu, Accent Conversion, Accent Classification, Machine Learning, Artificial Intelligence, Convolutional Neural Network, Communication barriers, Speech recognition, Text tospeech, Linguistic diversity, Communication technology, Natural language processing (NLP)

## I.INTRODUCTION

Our research delves into the intricate domain of "Telugu Accent Conversion and Classification Using AI." Focusing on the widely spoken Telugu language in Telangana and Andhra Pradesh, our study addresses the challenges posed by diverse accents within this linguistic landscape. Leveraging machine learning techniques, Convolutional Neural Networks (CNN) for accent classification and Google Textto-Speech (gTTS) for accent conversion, our approach aims to bridge the gap between diverse linguistic variations. The significance of this work lies not only in enhancing speech recognition systems but also in its potential application to other regional languages, offering a promising avenue for cross-cultural communication in the evolving landscape of AI-driven language processing.

## II.LITERATURE REVIEW

The research focused on investigating the effectiveness of diverse methodologies in a specific domain. Through examining relevant research papers, the goal was to assess various approaches and techniques employed in these areas. This process aimed to uncover nuanced intricacies and advancements within the field.

**Duduka et al. [1]** The study of Accent classification is a crucial aspect of language processing, contributing to the development of systems that can accurately identify and understand spoken accents. The utilization of neural networks in this context reflects a contemporary and best approach to solve the challenges in accent variability in speech. Neural networks, inspired by the structure of the human brain, excel at learning intricate patterns and representations from data, making them well-suited for complex tasks such as accent classification. the use of a neural network

suggests a data-driven approach. Neural networks are likely employed to extract relevant features from speech signals, enabling the model to discern subtle nuances associated with different accents. The choice of this approach indicates a departure from traditional methods and underscores the significance of leveraging advanced machine learning techniques in the field.The implications of this research extend beyond the academic realm, potentially influencing the development of real-world applications.

**Ensslin et al. [2]** The authors centered their research on the application of deep learning techniques for this purpose, opening up possibilities for enhancing the user experience in gaming environments through the adaptation to different accents .One of the main objectives of this research is likely to address the challenges associated with diverse accents encountered in online gaming scenarios. Accents can significantly vary based on players' geographic locations and linguistic backgrounds, leading to potential communication barriers and misunderstandings within multiplayer gaming environments. By leveraging deep learning, Ensslin and colleagues aimed to develop a system capable of accurately detecting and adapting to different accents in real-time. Deep learning, as a subset of machine learning, is well-suited for tasks involving complex patterns and hierarchical representations within data. In the context of speech accent detection, deep learning models, can learn intricate features and nuances associated with diverse accents, allowing for more accurate and robust classification. In online gaming, effective communication is often crucial for teamwork and collaboration. The adaptation might involve speech recognition systems adjusting their models to better understand and interpret the speech of players with various accents.

**Parikh et al. [3]** The primary aim of their work is likely centered around advancing systems capable of both recognizing and modifying accents in spoken English, with potential implications for a range of applications. Accent classification is a fundamental aspect of language processing, and the ability to accurately identify accents in spoken English is essential for various domains. It aids in tasks such as natural language understanding, and even language-based user interfaces. Parikh and their colleagues likely sought to leverage machine learning algorithms to enhance the precision and efficiency of accent classification within the English language. the mention of accent conversion in their work suggests a dual-purpose approach. In

addition to identifying accents, the research may involve developing techniques to modify or convert accents in spoken English. This could be of significance in applications where clear communication or standardized pronunciation is crucial, such as language learning platforms, voice assistants, or automated customer service systems. Machine learning, as applied in this research, is well-suited for handling the complexity of accent classification and conversion. By training models on a diverse dataset of English speech samples with various accents, the algorithms can learn patterns and features associated with different linguistic variations. The potential outcome of this research could be a system capable of automatically recognizing the accents of speakers and, if necessary, modifying the output to conform to a desired accent or pronunciation.

**Bird et al. [4]** The research, with a specific emphasis on both non-native and native speakers. The primary objective of their work appears to be the exploration of techniques for accurately classifying accents in the context of speech-based biometric systems, suggesting potential applications in security and authentication. Speech biometrics involves the use of distinctive voice characteristics for identification or verification purposes. Accent classification within this domain is a critical aspect, as accents can significantly influence the acoustic features of speech. Bird and their colleagues likely aimed to address the challenges associated with diverse accents, particularly in a security and authentication context where precision and reliability are paramount. Accurate accent classification in this context could have implications for the development of more inclusive and effective security systems that account for variations in speech patterns across different linguistic backgrounds.The potential applications of this research extend to security and authentication systems where voice-based identification is employed. By improving the accuracy of accent classification, the biometric systems can better adapt to the linguistic diversity of their user base. This could lead to enhanced security measures, reduced false positives or negatives, and a more user-friendly experience for individuals with different accents.

**Zhao et al. [5]** The research, showcasing a novel approach to transforming accents in speech signals. The central goal of their work appears to be the exploration of techniques that have the potential to facilitate cross-cultural communication and understanding by modifying the accents present in

spoken language. Accent conversion is a unique aspect of speech processing that goes beyond simple recognition or classification. Instead of merely identifying accents, Zhao and their colleagues delve into the transformation of spoken language to modify the accent characteristics. Phonetic posteriorgrams are likely employed as a representation of phonetic information, capturing the nuances of speech sounds in a way that is conducive to the accent conversion process. The application of accent conversion techniques aligns with the broader goal of improving communication across diverse linguistic backgrounds. In scenarios where individuals with different accents interact, misunderstandings may arise due to linguistic variations. Zhao et al.'s research aims to address this challenge by providing a method to convert spoken language to a more familiar or standardized accent, potentially fostering clearer and more effective communication. The use of phonetic posteriorgrams suggests a focus on the phonetic elements of speech, such as the pronunciation of specific sounds and the rhythm of spoken words. This level of granularity in representation is crucial for accurately capturing and modifying the phonetic characteristics associated with different accents.

**Gao et al. [6]** This research explores the application of GANs, a class of artificial intelligence models, to generate synthetic voices that convincingly mimic the characteristics of a target speaker.The use of GANs in voice impersonation signifies a departure from traditional approaches. GANs consist of two neural networks, a generator, and a discriminator, trained adversarially. The generator aims to create synthetic voices, while the discriminator seeks to distinguish between real and generated voices. This adversarial training process allows the model to iteratively refine its ability to produce authenticsounding impersonations, raising concerns about the potential for voice security breaches and unauthorized access. Voice impersonation techniques, if successful, could pose challenges to voice authentication systems, voice-controlled devices, and other applications relying on voice recognition.The specific methodologies employed by Gao et al. in their study, as well as the nuances of the GAN architecture used, are not explicitly detailed in the provided information.Voice impersonation techniques, particularly when facilitated by advanced machine learning models like GANs, could have far-reaching consequences in the realms of cybersecurity, privacy, and humancomputer interaction. Understanding the

capabilities and limitations of such technologies is crucial for ensuring the responsible development and deployment of voice-related applications.

**Desai et al. [7]** The research explores the application of artificial neural networks (ANNs) to voice conversion, a process aimed at transforming the characteristics of a speaker's speech to resemble that of a target speaker. Voice conversion is a crucial aspect of speech processing, with applications ranging from personalized voice synthesis to enhancing naturalness in voice-driven interfaces. The use of artificial neural networks in this context signifies a departure from traditional methods, leveraging the capacity of ANNs to learn complex mappings and representations from data.While the specific details of the neural network architecture and training methodologies employed in this study are not provided in the available information, the use of ANNs in voice conversion suggests a data-driven approach. ANNs are capable of capturing and learning the intricate features and patterns associated with individual voices, making them well-suited for tasks involving voice transformation.The implications of voice conversion using artificial neural networks extend to various domains. In personalized voice synthesis, individuals may have the opportunity to customize their synthesized voices for applications like voice assistants or navigation systems. Additionally, in voice-driven interfaces, achieving natural and contextually appropriate voice conversion is crucial for user acceptance and satisfaction.

**Aryal and Gutierrez-Osuna [8]** This research focuses on the application of artificial neural networks (ANNs) for accent conversion, a process aimed at modifying the accent characteristics in spoken language. The emphasis on artificial neural networks suggests a machine learning approach to address the challenges associated with accent conversion. Neural networks, particularly when applied to speech processing tasks, are adept at learning complex patterns and representations from data, making them suitable for tasks like accent modification. The technical report format suggests that the work may provide an in-depth exploration of the methodologies, experimental setup, and potentially the performance metrics used in the research. The implications of accent conversion using artificial neural networks are diverse, ranging from applications in language learning platforms to voice modification in communication systems. Such technology has the potential to facilitate cross-cultural communication and enhance the

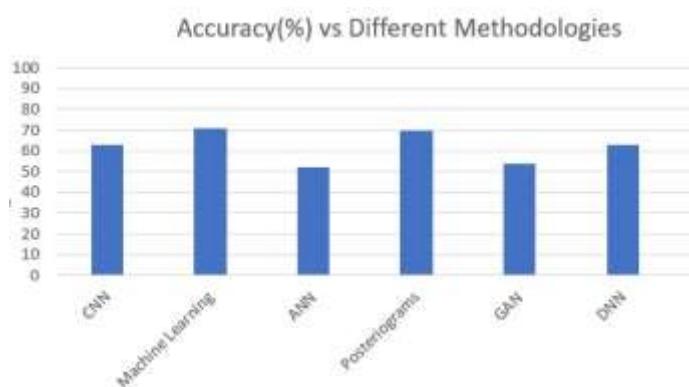adaptability of speech synthesis systems to different linguistic variations

**Badshah et al. [9]** This study focuses on the deep convolutional neural networks (CNNs) for the task of recognizing emotional states from spectrograms of speech signals. The use of deep CNNs in speech emotion recognition signifies an advancement in the field, leveraging the capability of these neural networks to automatically learn hierarchical features from input data. Spectrograms, which provide a visual representation of the frequency content of a speech signal over the time, serve as the input for the proposed emotion recognition system. While specific details regarding the architecture and training methodology of the deep CNNs are not available in the provided information, the utilization of CNNs in this context suggests an ability to capture intricate patterns in the spectrograms associated with different emotional states. CNNs are particularly effective in image-based tasks, and their application to spectrograms aligns with the visual nature of these representations. The implications of speech emotion recognition from spectrograms using deep CNNs are diverse. Emotion recognition can enhance human-computer interaction, sentiment analysis, and the development of emotionally intelligent systems. Applications include virtual assistants that respond appropriately to users' emotional states or systems that adapt content based on emotional feedback.

**Bearman, Josund, and Fiore [10]** This study delves into the use of deep neural networks (DNNs) for converting foreign accents, with a specific focus on articulatory-based methods.The emphasis on articulatory-based conversion suggests an approach that considers the physiological aspects of speech production, such as the movements of articulators . Deep neural networks are known for their ability to capture complex patterns and representations from data, and in this context, they are likely employed to learn the articulatory features associated with different accents.While the specific architectural details and training methodologies of the deep neural networks are not provided in the available information, the utilization of DNNs in articulatory-based conversion highlights the potential for these models to capture and synthesize the subtle nuances of foreign accents.The implications of articulatorybased conversion of foreign accents with DNNs extend to various applications. This approach may find utility in language learning platforms, cross-cultural communication tools, and systems that aim to adapt spoken language to specific linguistic variation.

| Paper | Year | Technology | Pros | Cons |
|---|---|---|---|---|
| 1 | 2021 | Advances accent classification using neural networks | Advances accent classification using neural networks | Complexity in implementation |
| 2 | 2017 | DL | Effective for speech detection in the video games | Application limited to gaming context |
| 3 | 2020 | Machine Learning | Addresses English accent classification and conversion | Complexity in implementation |
| 4 | 2019 | Accent Classification | Explores accent classification in human speech biometrics | - |
| 5 | 2018 | Phonetic Posteriorgrams | Introduces accent conversion using phonetic posteriorgrams | Complexity in implementation |
| 6 | 2018 | Generative Adversarial Networks | Addresses voice impersonation using generative adversarial networks | Ethical concerns, potential misuse |
| 7 | 2009 | Artificial Neural Networks | Investigates voice conversion using artificial | May require extensive training data |

| | | | neural networks | |
|---|---|---|---|---|
| 8 | 2015 | Deep Neural Network | Focuses on articulatory-based conversion of the accents | May require specialized data for non-native accents |
| 9 | 2017 | Deep Convolutional Neural Networks | Explores speech emotion recognition from spectrograms | May require substantial computational resources |
| 10 | 2017 | Artificial Neural Networks | May require substantial computational resources | Ethical concerns, potential misuse |


Accuracy(%) vs Different Methodologies

## III.CONCLUSION

In conclusion, the ongoing research addresses on Telugu accent conversion and classification using AI which marks a significant role towards overcoming communication challenges arising from diverse accents within the Telugu language. The integration of machine learning techniques, Convolutional Neural Network , has demonstrated promising results in accurately classifying various Telugu accents. This precision in accent classification lays a solid foundation for addressing the root cause of communication barriers among speakers with different linguistic backgrounds.

## REFERENCES

[1] Duduka S., Jain, H. Jain, H.P.V. and Chawan, P.M., 2021. A Neural Network Approach to AccentClassification. International Research Journal of Engineering and Technology (IRJET), 8(03), pp.11751177.

[2] Ensslin, A., Goorimoorthee, T., Carleton, S., Bulitko, V. and Hernandez, S.P., 2017, September.Deep Learning for Speech Accent Detection in Video Games. In Thirteenth Artificial Intelligence andInteractive Digital Entertainment Conference.

[3] Parikh, P., Velhal, K., Potdar, S., Sikligar, A. and Karani, R., 2020, May. English language accentclassification and conversion using machine learning. In Proceedings of the International Conferenceon Innovative Computing & Communications (ICICC).

[4] Bird, J.J., Wanner, E., Ekárt, A. and Faria, D.R., 2019, June. Accent classification in human speechbiometrics for native and non-native english speakers. In Proceedings of the 12th ACM InternationalConference on PErvasive Technologies Related to Assistive Environments (pp. 554-560).

[5]     Zhao, G., Sonsaat, S., Levis, J., Chukharev-Hudilainen, E. and Gutierrez-Osuna, R., 2018, April.Accent conversion using phonetic posteriorgrams. In 2018 IEEE International Conference onAcoustics, Speech and Signal Processing (ICASSP) (pp. 5314-5318). IEEE.

[6]     Gao, Y., Singh, R. and Raj, B., 2018, April. Voice impersonation using generative adversarialnetworks. In 2018 IEEE International Conference on Acoustics, Speech and Signal Processing(ICASSP) (pp. 2506-2510). IEEE.

[7]     Desai, S., Raghavendra, E.V., Yegnanarayana, B., Black, A.W. and Prahallad, K., 2009, April. Voiceconversion using artificial neural networks. In 2009 IEEE International Conference on Acoustics,Speech and Signal Processing (pp. 3893-3896). IEEE.

[8]     Aryal, S. and Gutierrez-Osuna, R., 2015. Articulatory-based conversion of foreign accents with deepneural networks. In the Sixteenth Annual Conference of the International Speech CommunicationAssociation.

[9]     Badshah, A.M., Ahmad, J., Rahim, N. and Baik, S.W., 2017, February. Speech emotion recognitionfrom spectrograms with deep convolutional neural networks. In 2017 international conference onplatform technology and service (PlatCon) (pp. 1-5). IEEE.

[10]    Bearman, A., Josund, K. and Fiore, G., 2017. Accent conversion using artificial neural networks.Stanford University, Tech. Rep, Tech. Rep.