Muskar Monish.et.al/ REST Journal on Data Analytics and Artificial Intelligence, 4(1), March 2025, 631-635



REST Journal on Data Analytics and Artificial Intelligence Vol: 4(1), March 2025 REST Publisher; ISSN: 2583-5564 Website: http://restpublisher.com/journals/jdaai/ DOI: https://doi.org/10.46632/jdaai/4/1/80



Cross-Language Speech Synthesis using Transfer Learning

* Muskar Monish, Meda Akhil, K. Uday, Vustela Pavani, S. Vijitha

Sunrise University, Alwar, Rajasthan, India. *Corresponding Author Email: muskarmonish@gmail.com

Abstract: The advancement of Text-to-Speech (TTS) technology has opened new possibilities for personalized voice applications, particularly in multilingual contexts. This project focuses on developing a cross-language speech synthesis system, converting English text into high-quality Telugu audio. The primary challenge lies in retaining the unique vocal characteristics of a specific speaker during this language transformation. By leveraging advanced deep learning models, such as Tacotron 2 and WaveGlow, along with transfer learning techniques, the system addresses linguistic and phonetic differences to generate natural and intelligible Telugu speech. This work bridges the gap between diverse languages in speech synthesis while contributing to advancements in personalized and multilingual voice applications.

Key Words: Cross-language speech synthesis, Deep Learning, Text-to-Speech (TTS), Transfer Learning, Speaker Identity, Taco Tron 2, Wave Glow.

1. INTRODUCTION

This project focuses on synthesizing Telugu speech from English text while retaining the speaker's vocal identity. Using advanced models like Tacotron 2 and WaveGlow, along with transfer learning, it addresses the challenges of cross-language speech synthesis, including differences in phonetics, prosody, and grammar. Tacotron 2, a sequence-to-sequence model, converts text into mel-spectrograms through its encoder-decoder structure, ensuring smooth and natural speech flow [1-4]. The model is pre-trained on English data and fine-tuned with Telugu datasets to adapt to the language's unique phonetics while preserving the speaker's voice characteristics. WaveGlow, a vocoder, transforms these mel-spectrograms into high-quality, natural-sounding waveforms with minimal computational overhead[5-8]. Voice cloning plays a crucial role in this project by using speaker embeddings to capture and retain the speaker's unique vocal traits such as tone, pitch, and style across both languages [11-13]. The training process begins with pre-trained English models, which are then fine-tuned with paired English text and Telugu audio to ensure accurate phoneme mapping and natural pronunciation [12][14]. The final system is evaluated based on clarity, naturalness, phonetic accuracy, and the preservation of speaker identity. The results demonstrate the success of this approach, achieving high-quality Telugu speech from English text while maintaining the speaker's identity. This innovation opens up possibilities for personalized, multilingual TTS systems [15], paving the way for advancements in real-time translation and voice assistant technologies

2. BACKGROUND

A.Cross-Language Text-to-Speech Synthesis: Cross-language TTS aims to synthesize speech in a language that differs from the speaker's native one, while maintaining their voice characteristics. This innovative approach requires not just understanding of multiple languages, but advanced modeling to preserve speaker identity. Recent advancements in deep learning have made significant strides in handling the challenges posed by linguistic variations, prosody, and phonetic differences across languages.

B. Deep Learning Models and Speaker Embedding: At the core of this technology are deep learning models like Tacotron and WaveGlow, which convert text to speech with high fidelity. Speaker embedding techniques are used to capture and replicate a speaker's unique characteristics, allowing for cross-language synthesis that retains the speaker's voice even when generating speech in a different language.

C. Challenges in Cross-Language TTS: One of the primary challenges lies in the phonetic discrepancies between languages, as TTS systems must handle different alphabets, tone, stress, and rhythm. Additionally, synthesizing speech from a language that a speaker has not recorded poses difficulties in maintaining naturalness and expressiveness. Addressing these issues requires large-scale datasets and advanced techniques in multilingual speech synthesis.

D. Applications and Impact: Cross-language TTS can revolutionize areas such as multilingual customer service, education, and accessibility. It allows a single speaker's voice to be used across different languages, enhancing the personalization of speech synthesis systems. The technology also promotes inclusivity by making digital interactions more accessible to speakers of various languages while preserving cultural and linguistic identity.

E. Future Challenges and Research Directions: Despite its potential, the field of cross-language TTS faces several challenges. These include ensuring high-quality synthesis for languages with limited data, fine-tuning models for accurate intonation, and addressing computational efficiency for real-time applications. Ongoing research focuses on improving these aspects and making the technology scalable to a wide range of languages and speakers [2]

3. LITERATURE REVIEW

Text-to-Speech (TTS) technology has made remarkable strides over the past few decades, with systems like Google's Tacotron 2, DeepMind's WaveNet, and Amazon Polly producing highly natural and humanlike speech. However, the majority of these systems remain primarily monolingual in design. Monolingual TTS systems are optimized to work in a single language, allowing for high-quality and natural-sounding speech. These systems have been highly effective for tasks such as virtual assistants, automated customer service, and accessibility tools for the visually impaired.

A.Tacotron 2: One of the most advanced neural TTS systems, Tacotron 2, generates speech by mapping sequences of characters directly to mel-spectrograms, which are then converted into speech using a vocoder like WaveGlow. Although Tacotron 2 is highly efficient at generating natural-sounding speech in English, it is limited by its monolingual focus. Its architecture is designed to optimize text-to-speech conversion for English, which has a distinct phonetic structure compared to other languages such as Telugu. The challenge with adapting this system for other languages lies in the phonetic, grammatical, and prosodic differences that vary significantly across languages. For example, while English uses stress and intonation in a relatively predictable way, Telugu has different patterns for stress and pitch, requiring the system to adapt to these nuances.

B.WaveNet: Another notable advancement in TTS is DeepMind's WaveNet, which generates speech by modeling audio at the waveform level. WaveNet can produce extremely natural-sounding speech by modeling the probability distribution of audio samples conditioned on the input text [16-18]. Like Tacotron 2, WaveNet has been designed for a single language, and while it excels in producing fluid and coherent speech, it requires large computational resources and large datasets for effective speech synthesis. When adapted to cross-lingual applications, WaveNet struggles to maintain vocal identity due to differences in language phonetics and grammar [19]. Moreover, training such a model for new languages is computationally intensive and requires substantial data collection efforts.

C.Amazon Polly: Amazon's Polly is a cloud-based service that converts text into lifelike speech. Polly is highly customizable, supporting a wide variety of languages and voices. However, it remains primarily monolingual in its architecture, requiring distinct models for each language. Polly's ability to switch between languages is still a challenge as it cannot seamlessly transition between languages while preserving speaker identity. Additionally, Polly's cross-lingual capabilities do not focus on retaining the vocal characteristics of the speaker, especially when generating speech in a language other than the original input.

The common thread among these systems is their focus on producing high-quality speech in a single language. Most TTS systems, including those mentioned, require vast amounts of training data to generate natural-sounding speech. They are trained on large datasets with thousands of hours of voice recordings, making them impractical for low-resource languages like Telugu, where such datasets may not be available [7][8]. Moreover, they are typically designed for specific language families, limiting their ability to transfer knowledge across vastly different linguistic structures.

While multilingual TTS systems have been developed, they often prioritize the generation of speech across different languages at the expense of maintaining a consistent vocal identity [20]. One area where these systems fall short is in cross-language speech synthesis, where the goal is to generate speech in a second language that retains the unique vocal traits of the original speaker. Existing systems like Tacotron 2, WaveNet, and Polly focus primarily on generating natural-sounding speech but do not prioritize maintaining speaker identity across languages [21]. The retention of speaker identity is crucial for applications that require personalized voice outputs, such as voice assistants, speech translation, and voiceovers. This presents a significant gap in current TTS technology, necessitating the development of new systems that can effectively handle these requirements.

4. FINDINGS AND LIMITATIONS

In this section, we will examine the results of our investigation in more detail. While existing TTS systems have achieved impressive results in generating lifelike speech, they possess several limitations when applied to cross-language tasks, especially when it comes to retaining speaker identity and adapting to low-resource languages.

- Monolingual Nature: Most TTS systems, including Tacotron 2, WaveNet, and Amazon Polly, are inherently monolingual. They are trained on large datasets specific to a single language and are optimized to perform well within that language's phonetic and grammatical framework. When applied to another language, the systems typically fail to capture the nuances of the target language, leading to unnatural-sounding speech. This limitation is particularly evident in cross-language scenarios, where the model's inability to transition smoothly between languages results in speech that does not sound coherent.
- Lack of Voice Consistency Across Languages: Maintaining voice consistency, especially across languages, is one of the most significant challenges in current TTS systems. While monolingual TTS models can replicate a speaker's voice with high accuracy in a single language, they often fail to retain the speaker's vocal identity when generating speech in another language. This is primarily because the models are not designed to preserve vocal characteristics (such as tone, pitch, and speaking style) across different phonetic structures. As a result, the synthesized voice in the second language may sound disjointed or different from the original speaker.
- High Resource Demand: Existing systems are highly resource-intensive. They require vast amounts of data to generate high quality speech, making them impractical for low-resource languages like Telugu. For instance, training a model like Tacotron 2 or WaveNet requires thousands of hours of high-quality voice recordings to produce convincing and fluent speech. Collecting such datasets is not always feasible, especially for less commonly spoken languages or for speakers with distinct accents. The computational power required to train these models is another barrier, as models like WaveNet, due to their fine-grained waveform generation, are particularly computationally expensive [3][4].
- Limited Support for Low-Resource Languages: Most existing systems focus on languages for which large, high-quality datasets are readily available, such as English, Mandarin, and Spanish. Low-resource languages like Telugu receive far less attention due to the difficulty in acquiring large datasets for training [9][11][13]. Even multilingual TTS systems, which aim to handle multiple languages, are often biased toward high-resource languages and tend to perform poorly on languages with limited data.

5. METHODOLOGY

The proposed methodology for "Cross-Language Speech Synthesis using Transfer Learning" involves leveraging state-of-the-art models, such as Tacotron 2 and WaveGlow, combined with transfer learning to generate high-quality Telugu speech from English text while preserving the speaker's unique vocal identity. The system is divided into several key components and follows a structured workflow to achieve its objectives [15-22]. The architecture of the system employs Tacotron 2 for generating mel-spectrograms from input text. Tacotron 2, a deep learning-based sequence-to-sequence model, consists of an encoder that processes input text into meaningful features and a decoder that generates mel-spectrograms frame by frame, ensuring smooth transitions and accurate prosody. To adapt this model for Telugu speech synthesis, transfer learning is applied. A pre-trained Tacotron 2 model, initially trained on large English datasets, is fine-tuned using a smaller Telugu dataset. This process enables the model to learn Telugu's unique phonetic and prosodic features while retaining the general knowledge of English speech synthesis. During fine-tuning, speaker embeddings are employed to preserve the unique characteristics of the speaker's voice, ensuring consistent vocal identity across languages. Once mel-spectrograms are generated, they are converted into waveforms

using WaveGlow, a flow-based vocoder. WaveGlow synthesizes high-quality, natural-sounding audio from the melspectrograms by directly modeling the distribution of waveforms. Fine-tuning WaveGlow on the same Telugu dataset ensures that the audio output captures Telugu's specific phonetic nuances, such as vowel harmony and retroflex consonants, while maintaining the naturalness of the speaker's voice. The combination of Tacotron 2 and WaveGlow ensures that the synthesized speech is clear, smooth, and free of artifacts, meeting the desired quality standards [23]. The implementation begins with data preparation, where English text is preprocessed by tokenizing it into characters or phonemes and aligning it with corresponding Telugu audio recordings. This step is crucial for ensuring accurate phoneme-to-speech mapping. During training, Tacotron 2 is fine-tuned on this paired dataset to adapt its internal parameters to Telugu phonology. Similarly, WaveGlow is fine-tuned to handle Telugu-specific spectrogram features and generate smooth waveforms. The system is designed to optimize performance with limited Telugu data, employing techniques like data augmentation and regularization to enhance model robustness. The methodology also includes rigorous testing and validation to ensure that the synthesized Telugu speech meets high standards of clarity, naturalness, and speaker identity retention. Subjective listening tests are conducted where evaluators assess the output for naturalness, phonetic accuracy, and identity preservation. Objective measures, such as speaker similarity scores and mel-cepstral distortion (MCD), are used to quantify the model's performance. Test cases involve both short phrases and longer paragraphs to evaluate the system's consistency and fluency across varying input complexities. By combining advanced deep learning models with transfer learning, the system successfully bridges the gap between two linguistically distinct languages, generating high-quality Telugu speech while retaining the original speaker's vocal traits. This methodology not only addresses the challenges of cross-language speech synthesis but also demonstrates the feasibility of developing multilingual text-to-speech systems for low-resource languages

6. FUTURE DIRECTION

Future directions for the cross-language speech synthesis system include expanding its capabilities to support multiple low-resource languages, such as Tamil, Kannada, and African or Indigenous languages, through fine-tuning and multilingual adaptations. Enhancing voice cloning by refining speaker embeddings to capture accents, dialects, and emotional tones, along with developing emotional speech synthesis, would improve realism and personalization. Optimizing the system for real-time applications, such as live translation tools and interactive virtual assistants, would enable instantaneous responses in multilingual scenarios. Integrating it into advanced conversational AI frameworks could enhance human-computer interactions, while reducing dependence on large datasets through techniques like unsupervised or zero-shot learning could enable scalability. Adding user customization features for voice tone, speed, and emotion would increase adaptability for diverse use cases, including accessibility and storytelling. Extensive user testing and evaluation in real-world scenarios, such as live translation or multilingual customer support, would further refine the system, ensuring it delivers high-quality, natural speech across various applications

7. CONCLUSION

The literature analysis concludes by highlighting the notable developments made in intrusion detection systems (IDS) for Industrial Internet of Things (IIoT) contexts. By combining advanced algorithms, hybrid models, and deep learning techniques, researchers have shown substantial improvements in efficiency and accuracy in detection. Furthermore, the emergence of specialized security models—like the proposed TSM—highlights the possibility of customized solutions to deal with certain IIoT security issues. However, issues including class imbalance, feature repetition, and dataset restrictions still exist and call for more research. The field of IIoT intrusion detection seems to have a bright future despite these obstacles, with chances to experiment with novel strategies and take advantage of cutting-edge technologies. Stakeholders may improve the security architecture of IIoT ecosystems by tackling these issues and adopting novel approaches to research.

REFERENCES

- Shen, J., Pang, R., Weiss, R. J., Schuster, M., Jaitly, N., Yang, Z., et al. (2018). Natural TTS Synthesis by Conditioning Wavenet on MEL Spectrogram Predictions. 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 4779–4783. https://doi.org/10.1109/ICASSP.2018.8461368
- [2]. Debnath, A., Patil, S. S., Nadiger, G., & Ganesan, R. A. (2020). Low-Resource End-to-end Sanskrit TTS using Tacotron2, WaveGlow and Transfer Learning. 2020 IEEE 17th India Council International Conference (INDICON), 1–5. https://doi.org/10.1109/INDICON49873.2020.9342071

- [3]. Tsai, W.-H., Thi, P. L., Tai, T.-C., Huang, C.-L., & Wang, J.-C. (2022). Low-Resource Speech Recognition Based on Transfer Learning. 2022 RIVF International Conference on Computing and Communication Technologies (RIVF), 145– 149. https://doi.org/10.1109/RIVF55975.2022.10013881
- [4]. Qiao, Z., Yang, J., & Wang, Z. (2023). Multi-Feature Cross-Lingual Transfer Learning Approach for Low-Resource Vietnamese Speech Synthesis. Proceedings of the 2023 3rd International Conference on
- [5]. Artificial Intelligence, Automation and Algorithms, 175–180. https://doi.org/10.1145/123456789.
- [6]. Weiss, R. J., Skerry-Ryan, R., Battenberg, E., Mariooryad, S., & Kingma, D. P. (2021). Wave-Tacotron: Spectrogram-Free End-to-End Text-to-Speech Synthesis. ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 5679–5683. https://doi.org/10.1109/ICASSP39728.2021.9413851
- [7]. Vinayasree, P., Reddy, A.M. A Reliable and Secure Permissioned Blockchain-Assisted Data Transfer Mechanism in Healthcare-Based Cyber-Physical Systems, Concurrency and Computation: Practice and Experience, 2025, 37(3), e8378
- [8]. Reddy, C.V.V., Reddy, A.M. An Optimized Transformative Approach To Detect Lung Cancer Based On Global & Local Feature Elements, Journal of Information Systems Engineering and Management, 2025, 10, pp. 494–505
- [9]. Vinayasree, P., Reddy, A.M. Deployment of Hybrid Chaotic Hashes for Blockchain Driven Internet 4.0 applications, Journal of Intelligent Systems and Internet of Things, 2025, 15(1), pp. 16–28
- [10].Samel, M., Mallikarjuna Reddy, A. An Enhanced Two-Way Unet Approach for Copy Move Image Forgery Detection, Journal of Information Systems Engineering and Management, 2025, 10(48), pp. 260–274.
- [11]. Vinayasree, P., Reddy, A.M. A Scalable And Secure Blockchain-Based Healthcare System: Optimizing Performance, Security, And Privacy With Adaptive Technologies, Journal of Theoretical and Applied Information Technology, 2024, 102(22), pp. 8084–8103
- [12].1. Gogu S, Sathe S (2022) autofpr: an efficient automatic approach for facial paralysis recognition using facial features. Int J Artif Intell Tools. https://doi.org/10.1142/S0218213023400055
- [13].2. Rao, N.K., and G. S. Reddy. "Discovery of Preliminary Centroids Using Improved K-Means Clustering Algorithm", International Journal of Computer Science and Information Technologies, Vol. 3 (3), 2012, 4558-4561.
- [14].3. Gogu, S. R., & Sathe, S. R. (2024). Ensemble stacking for grading facial paralysis through statistical analysis of facial features. Traitement du Signal, 41(2), 225–240.
- [15].Manoranjan Dash, N.D. Londhe, S. Ghosh, et al., "Hybrid Seeker Optimization Algorithm-based Accurate Image Clustering for Automatic Psoriasis Lesion Detection", Artificial Intelligence for Healthcare (Taylor & Francis), 2022, ISBN: 9781003241409
- [16]. Manoranjan Dash, Design of Finite Impulse Response Filters Using Evolutionary Techniques An Efficient Computation, ICTACT Journal on Communication Technology, March 2020, Volume: 11, Issue: 01
- [17].Manoranjan Dash, "Modified VGG-16 model for COVID-19 chest X-ray images: optimal binary severity assessment," International Journal of Data Mining and Bioinformatics, vol. 1, no. 1, Jan. 2025, doi: 10.1504/ijdmb.2025.10065665.
- [18].Manoranjan Dash et al.," Effective Automated Medical Image Segmentation Using Hybrid Computational Intelligence Technique", Blockchain and IoT Based Smart Healthcare Systems, Bentham Science Publishers, Pp. 174-182,2024
- [19].Manoranjan Dash et al.," Detection of Psychological Stability Status Using Machine Learning Algorithms", International Conference on Intelligent Systems and Machine Learning, Springer Nature Switzerland, Pp.44-51, 2022.
- [20].Samriya, J. K., Chakraborty, C., Sharma, A., Kumar, M., & Ramakuri, S. K. (2023). Adversarial ML-based secured cloud architecture for consumer Internet of Things of smart healthcare. IEEE Transactions on Consumer Electronics, 70(1), 2058-2065.
- [21].Ramakuri, S. K., Prasad, M., Sathiyanarayanan, M., Harika, K., Rohit, K., & Jaina, G. (2025). 6 Smart Paralysis. Smart Devices for Medical 4.0 Technologies, 112.
- [22].Kumar, R.S., Nalamachu, A., Burhan, S.W., Reddy, V.S. (2024). A Considerative Analysis of the Current Classification and Application Trends of Brain–Computer Interface. In: Kumar Jain, P., Nath Singh, Y., Gollapalli, R.P., Singh, S.P. (eds) Advances in Signal Processing and Communication Engineering. ICASPACE 2023. Lecture Notes in Electrical Engineering, vol 1157. Springer, Singapore. https://doi.org/10.1007/978-981-97-0562-7 46.
- [23].Ramakuri, S. K., Chithaluru, P., & Kumar, S. (2019). Eyeblink robot control using brain-computer interface for healthcare applications. International Journal of Mobile Devices, Wearable Technology, and Flexible Electronics (IJMDWTFE), 10(2), 38-50.
- [24].Daniel, G.V.; Chandrasekaran, K.; Meenakshi, V.; Paneer, P. Robust Graph Neural-Network-Based Encoder for Node and Edge Deep Anomaly Detection on Attributed Networks. Electronics 2023, 12, 1501. https://doi.org/10.3390/electronics12061501.