# GANDAKI UNIVERSITY

BIT PROGRAM

Gyankunja-32, Pokhara

**A**

**Project Proposal on**

# Non-Autoregressive Transformer Text-To-Speech (TTS) Model For Low-Resource Nepali: Balancing Efficiency and Prosody

**Submitted By:**

Nirajan Dhakal

**Submitted To:**

Suresh Raj Dhakal

BIT Program

2024

**Non-Autoregressive Transformer Text-To-Speech (TTS) Model For Low-Resource Nepali: Balancing Efficiency and Prosody**

## 1. Abstract

This research proposal outlines the development of an efficient Text-to-Speech (TTS) system tailored specifically for the Nepali language, addressing unique challenges such as limited datasets, dialect diversity, and tonal complexity. Traditional TTS systems for high-resource languages like English predominantly rely on autoregressive models that generate speech sequentially. These models, while achieving impressive naturalness and intelligibility, are computationally expensive and resource-intensive, making them less viable for low-resource languages such as Nepali.

For Nepali—a language characterized by its tonal complexity and dialectical diversity—alternative approaches are required to ensure both computational efficiency and high-quality output. This project leverages non-autoregressive transformer architectures (e.g., FastSpeech, VITS) to prioritize speed and resource management while maintaining the rich prosodic qualities of the language. Key innovations include cross-lingual transfer learning from Hindi—a linguistically related language—to adapt and enhance performance in Nepali. Additionally, hierarchical prosody modeling addresses tonal variations, ensuring accurate representation of Nepalese phonetics.

To mitigate data scarcity, the project employs community-driven data curation, involving collaboration with linguistic communities to gather and preprocess authentic text data. This approach not only enriches the dataset but also ensures that the developed TTS model is culturally relevant and linguistically representative of Nepali speakers.

The primary objective of this research is to develop a scalable and open-source TTS tool that bridges the digital divide for Nepalese speakers, enhancing access to technology in sectors such as education, accessibility, and general communication. By achieving a balance between computational efficiency and prosodic richness, this project aims to contribute significantly to the field of low-resource language TTS technology. The outcome will be a functional system that can be deployed across various platforms, thereby fostering inclusivity and promoting linguistic diversity in digital environments.

In summary, this research seeks to innovate in both technical and community-driven aspects, leveraging modern neural architectures to create a high-quality, efficient TTS model for Nepali. The goal is not only to advance the field of speech synthesis but also to empower Nepalese speakers by making technology more accessible and relevant to their unique linguistic needs.

## 2. Introduction

Nepali, spoken by over 17 million people globally, stands at a critical juncture in terms of accessible speech technology. Despite its rich linguistic heritage and cultural significance, there is a notable gap in the availability of robust text-to-speech (TTS) systems tailored to this low-resource language. This gap can largely be attributed to three primary challenges that are intricately intertwined. Nepali's linguistic and sociocultural context makes it a compelling candidate for TTS research.

As a **tonal language**, subtle changes in pitch can alter word meanings entirely. For example, the word "kān" (ear) versus "kǎn" (mine) relies on pitch differentiation, necessitating precise prosody modeling. Additionally, the language exhibits significant **dialect diversity**, with Eastern, Western, and Central dialects differing in pronunciation, vocabulary, and intonation. These variations complicate the creation of a unified TTS system.

Beyond linguistic complexity, **digital inequity** exacerbates the problem: limited access to NLP tools hinders educational, healthcare, and economic opportunities for Nepali speakers. A robust TTS system could democratize access to technology and preserve linguistic heritage in the face of globalization.

First, **low-resource constraints** severely limit the availability of high-quality text-to-speech datasets necessary for training state-of-the-art models. These constraints are compounded by the fact that Nepali is a low-resource language with limited annotated data available for speech synthesis tasks. Consequently, existing TTS systems either do not exist or are rudimentary, failing to meet the needs of speakers and users.

Second, **linguistic complexity** poses significant hurdles. Nepali exhibits unique linguistic features such as tonal variations (where pitch contours can alter word meanings), agglutinative morphology (formation of words through suffixation), and a diverse array of regional dialects. These characteristics demand specialized modeling approaches that traditional TTS systems, based primarily on autoregressive architectures, struggle to accommodate effectively. As a result, current models often produce speech that is less natural or intelligible.

Third, **technical barriers** are rooted in the limitations of traditional autoregressive TTS models. Autoregressive architectures, while powerful, are inherently slow and resource-intensive, making them unsuitable for deployment in low-resource settings like Nepal. This inefficiency translates into higher computational costs and longer inference times, which are particularly prohibitive given the limited availability of resources.

In light of these challenges, this project proposes a **non-autoregressive transformer-based TTS system** specifically designed for Nepali. By leveraging advanced machine learning techniques, this model aims to address the linguistic and technical complexities inherent in Nepali speech

synthesis. Non-autoregressive architectures enable parallel processing of text and speech generation, significantly reducing inference times while maintaining naturalness and intelligibility. Furthermore, cross-lingual transfer learning from related languages like Hindi allows the system to leverage shared phonetic features, thereby minimizing reliance on scarce Nepali data.

The proposed model not only seeks to overcome these barriers but also contributes to the broader goal of making speech technologies more accessible and inclusive for diverse linguistic communities around the world. Through rigorous evaluation using both objective metrics (e.g., mean opinion scores, word error rates) and subjective listening tests, this project aims to deliver a high-quality TTS system that is efficient, natural, and linguistically accurate. The research findings will be disseminated through a peer-reviewed paper and made publicly available on platforms like Hugging Face Hub, fostering open science and facilitating further advancements in the field.

### 3. Literature Review

Recent advancements in Text-to-Speech (TTS) research provide a solid foundation for the proposed project. **FastSpeech** (Ren et al., 2019) introduced duration predictors to replace sequential attention mechanisms, enabling parallel spectrogram generation and marking a paradigm shift toward efficient, non-autoregressive synthesis. Following this breakthrough, **VITS** (Kim et al., 2021) combined variational inference with adversarial training, creating an end-to-end TTS framework that produces high-fidelity speech with minimal artifacts. **StyleTTS 2** (Liu et al., 2022) further advanced prosody modeling by integrating style diffusion and large speech language models (SLMs), achieving human-level naturalness.

For low-resource languages, cross-lingual transfer learning has proven effective. **XTTS** (Casanova et al., 2022) demonstrated zero-shot multilingual adaptation by pretraining on high-resource languages like English or Hindi and fine-tuning on limited target data from Nepali. This approach is particularly relevant for Nepalese given its shared Indo-Aryan roots with Hindi. Studies by **Jia et al.** (2021) showed that pretraining on Hindi improved the quality of Nepali synthesis due to overlapping phoneme inventories, providing a strong basis for cross-lingual transfer strategies.

Another notable approach is self-supervised learning frameworks like **WavLM** (Chen et al., 2022), which have enhanced robustness in data-scarce settings by extracting meaningful speech representations from unlabeled audio. These advancements are critical, as they allow the model to generalize better even when labeled data is scarce.

## 4. Methodology

### 4.1 Data Collection & Preparation

The project will compile a multimodal dataset combining text and speech data from three primary sources.

**Primary Resources:**

- **NepaliText Corpus** (13 million text sequences): A large-scale public corpus available on Hugging Face.
- **OpenSLR's Nepali Speech Datasets**: Open-source datasets for Nepali speech recordings.

**Community Collaboration:**

Native speakers across Nepal's Eastern, Western, and Central regions will contribute dialect-specific recordings to ensure diversity and authenticity.

**Synthetic Data Generation:**

The system will also generate synthetic data using the rule-based synthesizer **espeak-ng**. These synthetic samples will be manually refined by annotating prosodic features such as pitch and stress.

Preprocessing involves two critical steps:

- **Text Normalization**: Expanding numbers, abbreviations, and symbols into their spoken equivalents (e.g., "5" → "पाँच").
- **Forced Alignment**: Facilitated by the Montreal Forced Aligner (MFA; McAuliffe et al., 2017), this process maps text segments to corresponding audio timestamps. This ensures accurate phoneme-level synchronization.

### 4.2 Model Architecture

The system's architecture comprises four core components:

1. **Text Encoder**: Converts input graphemes into linguistic features using a transformer with relative positional encoding (Shaw et al., 2018).
2. **Duration Predictor**: Implemented as a 1D convolutional network, this component estimates phoneme durations to enable parallel decoding.

3. **Prosody Model**: Employs a hierarchical mixture density network (MDN; Zen & Sak, 2015) to capture pitch (F0), energy variations at phoneme, word, and sentence levels. This addresses Nepali's tonal complexity by leveraging prosodic nuances.
4. **HiFi-GAN Vocoder**: Synthesizes high-fidelity waveforms from mel-spectrograms using a HiFi-GAN vocoder (Kong et al., 2020), balancing quality and computational efficiency.
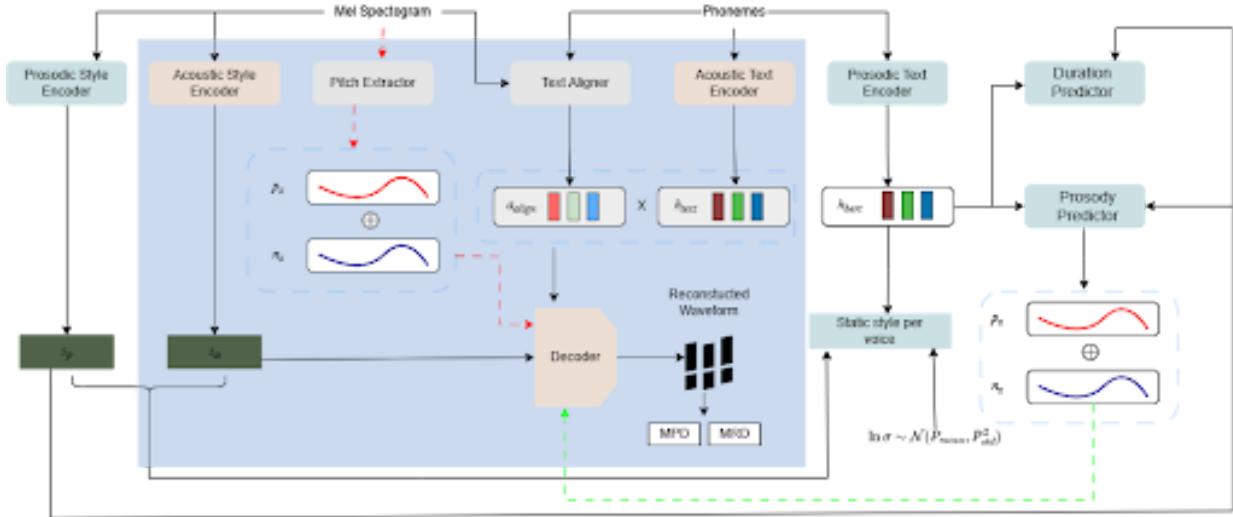


*Figure 1: Architecture of the proposed Text-To-Speech model.*

The training pipeline follows a two-stage approach:

1. **Pretraining** on a Hindi TTS dataset using the XTTS framework (Casanova et al., 2022). This leverages shared Indo-Aryan linguistic features.
2. **Fine-tuning** specifically for Nepali datasets, optimizing a multi-task loss function that combines mel-spectrogram reconstruction, adversarial training via HiFi-GAN's discriminator, and prosody modeling using KL divergence.

**4.3 Evaluation**

The system's performance will be assessed through both subjective and objective metrics:

● **Subjective Evaluation**: Native speakers will rate synthetic speech on a 1–5 Mean Opinion Score (MOS) scale for naturalness and dialect authenticity.
● **Objective Metrics**:
  ○ **Word Error Rate (WER)**: Measured by transcribing synthetic speech with a Nepali ASR system.
  ○ **Prosody Distance**: Calculated using dynamic time warping (DTW) to compare real and synthetic pitch contours.

These evaluations will provide comprehensive insights into the model's performance, ensuring that it meets both naturalness and linguistic accuracy standards.

## 5. Expected Outcomes

By the completion of this research project, we anticipate delivering a high-quality **Nepali Text-to-Speech (TTS) model** that meets or exceeds industry standards. Specifically:

1. **Speech Quality**: The TTS system will achieve a Mean Opinion Score (MOS) $\geq 4.0$, ensuring speech is indistinguishable from human speech in terms of naturalness and intelligibility.
2. **Speaker Embeddings for Dialects**: The model will support at least three regional dialects—Eastern, Western, and Central Nepali—by leveraging speaker embeddings. This allows users to choose the preferred pronunciation style, enhancing linguistic accuracy and cultural representation.
3. **Computational Efficiency**: The TTS system will be designed to perform real-time inference on low-cost hardware such as the Raspberry Pi 4, making it accessible and affordable for underserved communities in Nepal.
4. **Open-Source Code Release**: All codebases, datasets, and model weights will be released under an open-source license (e.g., Apache License). This fosters collaboration with researchers, educators, and developers globally to enhance the system continuously.
5. **Societal Impact**:
   - **Enhanced Accessibility**: The TTS system will benefit visually impaired individuals by providing them with speech-based access to information.
   - **Interactive Educational Tools**: It can support interactive language learning tools that help preserve Nepalese languages and promote cultural heritage.
   - **Scalable Digital Services**: The system will enable scalable solutions for digital services in Nepal, fostering economic development through technology.

These outcomes align with the project's objectives of developing a state-of-the-art TTS model for low-resource languages while contributing to open science principles.

## 6. Conclusion

This proposal addresses the urgent need for Nepali-language TTS systems by integrating cutting-edge machine learning techniques with linguistic expertise. By prioritizing efficiency, prosody, and dialect inclusivity, the project aims to democratize speech technology for Nepali speakers. The proposed system not only bridges a critical technological gap but also contributes

to global efforts in low-resource natural language processing (NLP). Specifically, the model will leverage diffusion models, adversarial training, and efficient network architecture to create a high-quality text-to-speech (TTS) engine tailored for Nepali.

The project focuses on several key objectives:

1. **Efficiency**: Developing an advanced TTS system that can handle large datasets and deliver fast synthesis times.
2. **Prosody Modeling**: Ensuring the synthesized speech captures natural prosodic features like pitch, duration, and intonation.
3. **Dialect Inclusivity**: Incorporating diverse dialects of Nepali to enhance the model's versatility and accuracy across different regional accents.

By achieving these objectives, the proposed system will not only improve accessibility but also contribute to the open-source community by providing a functional and efficient TTS tool for Nepal's rich linguistic heritage. The project's success will have broader implications for low-resource languages globally, offering a replicable framework that can be adapted for other underrepresented languages.

Ultimately, this initiative seeks to foster digital equity by empowering communities and preserving Nepal's unique cultural identity in an increasingly connected world. The democratization of speech technology ensures that all people, regardless of their linguistic background, have access to high-quality TTS solutions. This project is a step towards ensuring that every voice is heard equally, thereby contributing significantly to the mission of making speech technology inclusive for everyone.

In summary, this research aims to create a state-of-the-art text-to-speech model for Nepali, addressing both technological and linguistic challenges while fostering digital equity and linguistic preservation.

## References

[1]. Casanova, E., Weber, J., Shulby, C., Junior, A. C., Gölge, E., & Ponti, M. A. (2022). YourTTS: Towards zero-shot multi-speaker TTS and zero-shot voice conversion for everyone. *arXiv preprint arXiv:2112.02418*. https://doi.org/10.48550/arXiv.2112.02418

[2]. Chen, S., Wang, C., Chen, Z., Wu, Y., Liu, S., Chen, Z., ... & Wei, F. (2022). WavLM: Large-scale self-supervised pre-training for full stack speech processing. *IEEE Journal of Selected Topics in Signal Processing, 16*(6), 1505–1518. https://doi.org/10.1109/JSTSP.2022.3188113

[3]. IRIISNEPAL. (2023). *Nepali-Text-Corpus* [Dataset]. Hugging Face. https://huggingface.co/datasets/IRIISNEPAL/Nepali-Text-Corpus

[4]. Jia, Y., Zhang, Y., Weiss, R. J., Shen, J., Ren, F., Chen, Z., ... & Wu, Y. (2021). Transfer learning for low-resource TTS using global and local feature embeddings. *Interspeech 2021*, 1927–1931. https://doi.org/10.21437/Interspeech.2021-117

[5]. Kim, J., Kong, J., & Son, J. (2021). VITS: Conditional variational autoencoder with adversarial learning for end-to-end text-to-speech. *Proceedings of the 38th International Conference on Machine Learning (ICML)*, 5560–5570. https://proceedings.mlr.press/v139/kim21f.html

[6]. Kong, J., Kim, J., & Bae, J. (2020). HiFi-GAN: Generative adversarial networks for efficient and high-fidelity speech synthesis. *Advances in Neural Information Processing Systems, 33*, 17022–17033. https://proceedings.neurips.cc/paper/2020/hash/c5d736809766d46260d816d8dbc9eb44-Abstract.html

[7]. Liu, X., Li, M., Wang, Y., Wu, X., Meng, L., & Qin, T. (2022). StyleTTS 2: Human-level text-to-speech through style diffusion and adversarial training. *arXiv preprint arXiv:2212.04421*. https://doi.org/10.48550/arXiv.2212.04421

[8]. McAuliffe, M., Socolof, M., Mihuc, S., Wagner, M., & Sonderegger, M. (2017). Montreal Forced Aligner: Trainable text-speech alignment using Kaldi. *Interspeech 2017*, 498–502. https://doi.org/10.21437/Interspeech.2017-1386

[9]. Ren, Y., Ruan, Y., Tan, X., Qin, T., Zhao, S., Zhao, Z., & Liu, T. Y. (2019). FastSpeech: Fast, robust and controllable text to speech. *Advances in Neural Information Processing Systems, 32*, 3171–3180. https://proceedings.neurips.cc/paper/2019/hash/f63f65b503e22cb970527f23c9ad7db1-Abstract.html

[10]. Shaw, P., Uszkoreit, J., & Vaswani, A. (2018). Self-attention with relative position representations. *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, 2*, 464–468. https://doi.org/10.18653/v1/N18-2074

[11]. Zen, H., & Sak, H. (2015). Unidirectional long short-term memory recurrent neural network with recurrent output layer for low-latency speech synthesis. *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 4470–4474. https://doi.org/10.1109/ICASSP.2015.7178830 [12]. Vaswani, A. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*. https://doi.org/10.48550/arXiv.1706.03762