

Prosody Predictor based Diffusion Models Techniques for Enhanced Speech Synthesis

Dr. K. Aruna Bhaskar^{1,*}, Dr. Bechoo Lal², Dr. M Bhaskar³, S. Sushma⁴, N. Praveen⁵, A. Siva Kumar Reddy⁶

^{1,2,3,6}Department of Computer Science and Engineering, Koneru Lakshmaiah Education Foundation, Vaddeswaram, Guntur Dist., Andhra Pradesh - 522302, India.

⁴Department of Information Technology, Aditya University, Surampalem, Andhra Pradesh - 533437, India.

⁵Department of Computer Science and Design, Sagi Rama Krishnam Raju Engineering College, Bhimavaram, Andhra Pradesh, India

¹arunabhaskar@kluniversity.in, ²bechoolal@kluniversity.in, ³bhaskarmarapelli@gmail.com, ⁴sushma.cse2@gmail.com,

⁵praveen@srkrec.edu.in, ⁶skumar_a007@yahoo.com

ARTICLE INFO

ABSTRACT

Received: 26 Dec 2024

Revised: 14 Feb 2025

Accepted: 22 Feb 2025

A prosody predictor based on a diffusion model is crucial to the new zero shot approach of voice synthesis. Since diffusion models excel at capturing complicated distributions, they are perfect for simulating the complex patterns of prosody in speech. These models have recently attracted interest in a number of generative tasks. By repeatedly changing an initial chaotic input into an output that nearly matches the intended goal, a diffusion model acts by gradually refining the input. The diffusion model iteratively refines an initial rough estimate of the prosody pattern in the context of prosody prediction. To get realistic sounding speech, it is necessary to capture small prosody fluctuations in pitch, length, and loudness. This approach enables the model to do just that. Training on massive speech corpora teaches the diffusion model-based prosody predictor to mimic reference speech in its prosody pattern generation. During inference, the model makes use of the learnt prosody patterns to anticipate the target speech's prosody, guaranteeing that the produced speech is expressive and authentic, even while the speaker is unseen.

Keywords: Text to Speech (TTS), Machine Learning, Artificial Intelligent (AI), Prosody Model, Diffusion Model.

INTRODUCTION

One goal of text to speech (TTS) synthesis is to make the translated text sound as human like and natural as feasible. Making speech that sounded both natural and comprehensible was a problem for early TTS systems. Nevertheless, these systems have been greatly enhanced the advent of neural networks and deep learning has empowered them to generate speech which is almost comprehensible and comparable in naturalness to human speech. In the training dataset, where these innovations have been most noticeable, the generated speech for speakers is almost inaudible compared to the natural human voice. There are still several limits to existing TTS systems, even with these improvements [1]. The fact that these systems are often limited to the speaker's that were part of the training data is a major drawback.

These audio recordings are frequently recorded in controlled environments, such as professional recording studios, to ensure exceptional audio quality. The issue is that this approach restricts the system's ability to replicate timbre, prosody, and modality-focused analysis of human speech [2]. To put it simply, text-to-speech algorithms are very good at mimicking the natural speech of speakers they have already been trained on, but they struggle to do the same for unknown or unfamiliar speakers [3]. This limitation significantly reduces the adaptability and general use of a TTS system in real-world scenarios where it could be necessary for the system to generate speech for a range of voices. Zero shot speech synthesis has consequently gained popularity as a topic of discussion. A text's ability to speak (TTS). Zero shot speech synthesis is the ability of a text-to-speech (TTS) system to generate speech for any speaker, even if the speaker was not trained with the system's training data [4].

The new speaker's acoustic reference is all we want to use for this. This enables the system to use a brief audio sample of the new speaker's voice to replicate their speaking pattern. Although there are a number of challenges with zero shot speech synthesis, this feature is highly desired because it would make TTS systems much more flexible and adaptive, allowing them to produce speech that accurately represents any speaker, regardless of whether their voice was present in the original training data or not [5] [6]. Developing accurate models of speaker timbre—the unique quality or tone of a human voice—is the primary challenge. A person's voice might sound noticeably different depending on their timbre, even when they are using the same phrases. For zero shot speech synthesis, the system must be able to accurately capture and duplicate the speaker's timbre in order to ensure that the synthesized speech sounds like it truly comes from an unseen speaker [7].

This essay describes why recent advancements in zero shot text to speech (TTS) algorithms have garnered a lot of attention due to their ability to generate speech that sounds very expressive and natural. These methods increase the generalizability of TTS models by utilizing large datasets, which enables them to capture additional speaker characteristics such as prosody, style, and timbre. Because of this, the models are able to imitate the speech of unfamiliar speakers even when they were not exposed to their speech during training [8].

Discrete token-based and continuous vector-based zero shot TTS algorithms are the two primary types of these algorithms. What sets these approaches apart is how they analyze and portray the acoustic data. VALL E is a well-known method that makes use of discrete tokens. VALL E uses a neural code that uses Residual Vector Quantization (RVQ) to tokenize spoken language. This method autoregressively generates the target speech tokens from both text and audio tokens using a language model [9]. The autoregressive method allows speech to be produced sequentially by leveraging the outcomes of previous token forecasts.

This approach has demonstrated impressive results in capturing speaker similarity for speakers that were not included in the training set. On the other hand, SPEAR TTS is another zero shot TTS method that uses the existing Audio LM framework to provide a gradual text-to-speech transition [10]. In short, SPEAR TTS converts text into acoustic tokens that represent speech sounds and then back into semantic tokens that represent the text's meaning. Even when the reference speaker is not present in the training data, this method aims to approximate the reference speaker's sound as much as feasible, much like VALL E [11].

has significantly influenced advancements in deep learning applications for speech processing as well as neural network-based voice synthesis and recognition systems [12]. A top-K recommender system that effectively handles

missing data is built using an ensemble technique. Conventional recommender systems, like collaborative filtering and content-based techniques, can suffer from sparse or missing user-item interaction data. In order to address this challenge [13]. A system for text-to-speech (TTS) synthesis designed specifically for Kannada, a Dravidian language spoken mostly in Karnataka, India. The project aims to improve Kannada speech synthesis by developing a linguistically and phonetically correct TTS system that generates a natural-sounding voice from text input [14]. MFCCs, which are crucial components of speech processing, represent the frequent characteristics of human speech. The study proposes a neural network-based approach to improve the accuracy and naturalness of Arabic speech synthesis [15].

FastSpeech is a new non-autoregressive text-to-speech (TTS) model designed to increase the speed, robustness, and controllability of speech synthesis. Traditional TTS models, such as Tacotron, use autoregressive processes that generate speech sequentially. These processes have several shortcomings, such as slow inference speeds, errors like mispronunciations, and a lack of precise control over prosody. [16] [17]. Historically, SPSS has struggled to generate natural-sounding speech because of over-smoothing effects. By incorporating GANs, the authors improve the expressiveness and naturalness of synthetic speech, producing a more realistic output waveform [18]. An end-to-end text-to-speech (TTS) system called NaturalSpeech seeks to provide voice synthesis that is as good as that of a human. The system generates expressive, high-quality speech from text input by utilizing deep learning innovations such as transformers and neural network-based models. The goal of NaturalSpeech is to bridge the gap between machine-generated speech and natural human speech by improving the expressiveness, naturalness, and intelligibility of synthesized speech [19]. text-to-speech (TTS) conversion systems, highlighting the basic concepts and challenges involved in developing effective TTS technology [20].

The TD-PSOLA technique is used to synthesis speech from factors like pitch and duration and allows for efficient compression of speech signals. The authors look at how the TD-PSOLA technique may be used to low bit rate speech coding systems, which makes it a great tool for communication systems with limited bandwidth [21]. By generating remarkably realistic and high-quality audio directly from raw audio waveforms, WaveNet provides significant advancements over traditional speech synthesis methods like parametric synthesis or concatenative synthesis [22]. To improve the voicing, naturalness, and intelligibility of whispered speech, generative models are employed. Whispered speech can sound odd or be difficult to understand because it lacks the tone and resonance of normal speech, which makes it challenging to synthesize [23].

A popular architecture in natural language processing (NLP), the Seq2Seq (Sequence-to-Sequence) model serves as the basis for the Deep Text-to-Speech (TTS) system [24]. employing a quantized Fo modeling approach for text-to-speech synthesis that is based on a Recurrent Neural Network (RNN). Since it influences the speech's pitch, the fundamental frequency, or Fo, is an essential part of speech and is necessary to give artificial speech a more expressive and natural tone [25]. Training models often requires large amounts of labeled data, which can be costly and time-consuming to get. One of the challenges in creating TTS is this. The authors offer methods for developing TTS systems using publicly available speech datasets or text corpora [26].

comprehensive examination of the state-of-the-art in speech synthesis, encompassing the various methods used to generate natural-sounding speech from text [27]. the creation of a Malayalam text-to-speech (TTS) system for Android that makes use of the concatenative synthesis technique. Due to its complex phonetic and syntactic patterns,

Malayalam, a language spoken mostly in the Indian state of Kerala, is a challenging language to create a high-quality TTS system for [28]. For Malayalam, which is quite similar to Telugu, concatenative TTS is utilized, particularly in regards to the concatenative synthesis process [29]. Important information regarding the use of concatenative synthesis to the Dravidian language may be found in the Kannada language. One of the several Indian language TTS systems that emphasizes the importance of concatenative synthesis is Telugu. In Indian languages like Telugu, concatenative synthesis is used to highlight the challenges and solutions of creating speech that sounds natural. [30] [31].

Natural Speech 2 uses continuous vectors to represent data. Natural Speech 2 uses a neural codec architecture that is similar to VALL E, which is based on RVQ, when encoding audio into continuous vectors. These continuous vectors, which capture subtle variations in prosody and style, provide a more accurate representation of the speech data. To forecast these continuous vectors, the Natural Speech 2 system uses a latent diffusion model. Because this model generates the continuous vectors through a process of gradual refinement, the generated speech can accurately replicate the characteristics of the reference speaker.

OBJECTIVES

This work aims to create a prosody predictor for zero-shot voice synthesis based on a diffusion model. The diffusion model was selected because it can capture intricate prosody patterns, such as changes in loudness, duration, and pitch, all of which are necessary to produce speech that sounds natural. The model learns to produce expressive prosody patterns that closely mimic those in human speech by repeatedly improving a preliminary rough approximation. The model can generalize to unknown speakers thanks to training on large speech datasets, guaranteeing that the synthesized speech is expressive and authentic even in zero-shot situations.

METHODS

The model can capture complex prosodic patterns, leading to more expressive and lifelike speech, by combining a hierarchical prosody adaptor with a diffusion model-based prosody predictor. **Accurate Speaker Timbre Reproduction:** Even in situations when the reference voice is absent from the training set, the model can faithfully mimic the timbre of speaker's that have not been observed by integrating a speaker encoder that generates a global speaker embedding.

1. **Less Dependency on Quantization:** The model steers clear of over quantizing speech, which may cause significant timbre features to be lost. Rather, it uses the global speaker embedding as a guide to produce speech that is timbrally correct.
2. **Scalability to New Speaker's:** The model is well suited for zero shot situations, in which the objective is to produce speech for new speakers with the least amount of training data possible. This is because of its capacity to generalize across various voices.
3. **Enhanced Naturalness and Expressiveness:** The model may produce speech that is almost identical to the reference speaker's naturalness and expressiveness by modelling prosody hierarchically and honing it via diffusion.

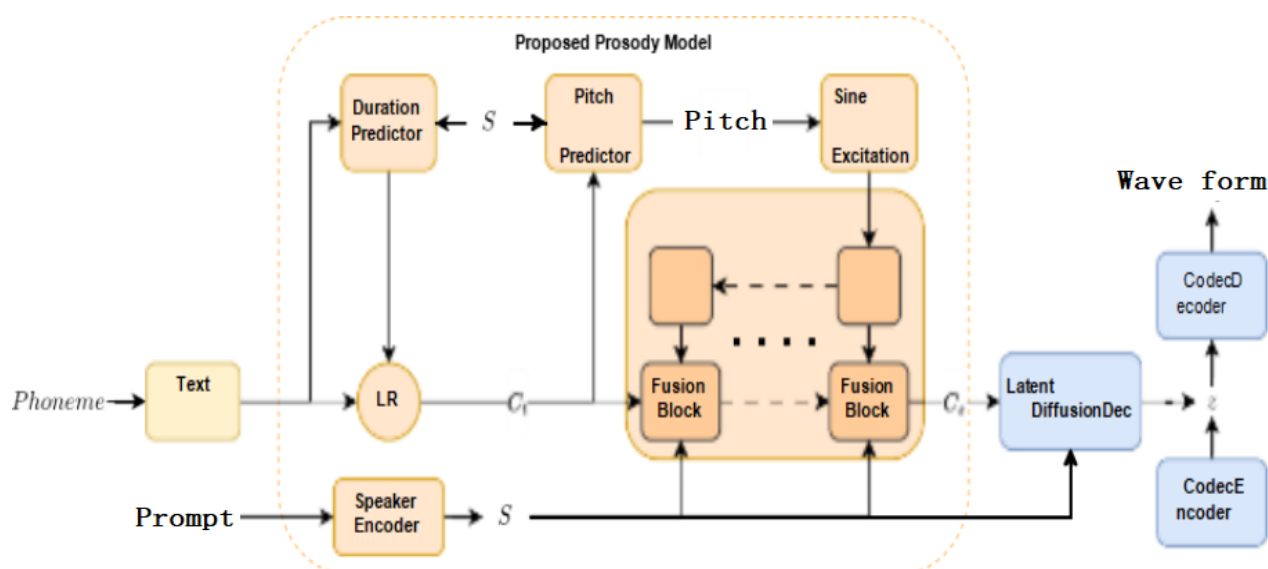


Figure 1. Enhanced Prosody Modelling

The proposed zero shot speech synthesis model addresses the shortcomings of current techniques in speaker timbre reproduction and prosody modelling, which is a major improvement in the area of TTS. The model is able to capture complex prosody patterns via the use of a hierarchical prosody adaptor and a diffusion model-based prosody predictor. Additionally, the speaker encoder guarantees that the produced speech retains the reference speaker's timbre. Through the use of extensive corpora and a decreased dependence on quantized audio, the model is capable of producing expressive and genuine speech, even for speaker's it has not heard in training. This method opens the door for more precise and adaptable text to speech systems that may be used for a variety of tasks, such as content generation and virtual assistants.

Prosody Predictor Based Diffusion Models Techniques

Machine learning has come a long way since diffusion models were first used, especially in generative tasks where producing a wide variety of high-quality outputs is of the utmost importance. The difficult job of forecasting pitch a crucial component of speech prosody is one area where diffusion models have shown remarkable potential, and text to speech (TTS) synthesis is one such domain. A speaker's pitch the perceived frequency of their speech sounds is crucial for expressing tone, meaning, and authenticity in their speech. To create expressive and natural sounding synthesized speech, it is crucial to capture the minute changes in pitch that happen in human speech. This has shown to be a particularly challenging problem for traditional methods of pitch prediction in TTS, particularly when confronted with complex and varied prosodic patterns. To improve the variety and authenticity of pitch predictions in voice synthesis, the diffusion model and more specifically, the denoising diffusion probabilistic model (DDPM) becomes useful.

The capacity of a diffusion model to train a strong generative process that can capture the complex subtleties of pitch changes is a major benefit when utilizing it for pitch prediction. Deterministic methods, in which one fixed pitch contour is produced for each input, are often used by traditional models. One consequence of this is the loss of expressiveness in speech due to over smoothing, which occurs when the natural fluctuation in pitch is artificially altered. Diffusion models, on the other hand, are probabilistic, so they can test out more variations in pitch, which

ultimately results in outputs that are more accurate representations of the inherent variety in human speech. The diffusion model is trained to anticipate the reversal of noise addition, which allows it to learn to provide a range of potential pitch contours instead of a single predictable one. The end product is synthetic speech that is more alive, interesting, and realistic in its prosodic patterns.

1. The diffusion model's capacity to manage the hierarchical and multi scale aspects of prosody is an additional major advantage when it comes to using it for pitch prediction. There are several granularities of pitch variation, ranging from very small changes within phonemes to larger patterns of intonation that cover whole phrases. Because it generates data iteratively and step by step, the diffusion model is ideal for modelling these varying prosody levels.
2. The model is able to capture both coarse grained and fine-grained fluctuations in pitch since it adds another layer of information to the pitch prediction with each phase of the reverse process. With this multi scale method, the prosodic structure of the synthesized speech is preserved on a global scale while the local pitch patterns are accurately captured. It is crucial for zero shot TTS situations to have this hierarchical modelling capacity since the model has to be able to be generalized to new speaker's and linguistic settings without previous exposure. The diffusion model's resilience in dealing with various emotional expressions and speaking styles is enhanced by its capacity to produce a wide range of pitch patterns. Pitch is a powerful tool for conveying meaning and emotion in human speech.
3. A speaker may, for instance, show enthusiasm by raising their pitch or gravity by lowering it. In order to create expressive and emotionally resonant synthesized speech, it is essential to capture these variances. Because of its probabilistic character, the diffusion model can manage the pitch variability associated with various emotional states and speaking styles better by exploring a broad variety of pitch shapes. This improves the TTS systems expressiveness by creating synthetic speech that sounds natural and communicates the speaker's desired emotional content.

In addition, the generating process of the diffusion model enables it to successfully manage the inherent variability and uncertainty in pitch prediction. Many things influence the precise pitch contour of a phrase in natural speech, including as the speaker's emotional state, the topic at hand, and the weight assigned to individual words. This unpredictability is a common challenge for traditional pitch prediction methods, which may result in forecasts that are either excessively stiff or too smooth. Instead of producing a single, fixed result, the diffusion model is structured to deal with this uncertainty by producing a distribution of potential pitch contours. More natural sounding speech is the result of the model's ability to generate realistic and diverse pitch patterns. There are a number of real-world benefits to training and deploying a diffusion model for pitch prediction. The model may be trained to provide high-quality outputs with a reduced amount of training data attributable to the iterative nature of the diffusion process.

When dealing with zero shot TTS, this becomes even more crucial since the model has to be able to generalize to different speakers and situations with less training data. For TTS systems that need to work effectively in a variety of situations, the diffusion model is a great tool since it can learn a robust generating process with little input. Furthermore, the diffusion model is a strong and adaptable tool for boosting the expressiveness and naturalness of synthesized speech. Its flexibility makes it easy to incorporate with current TTS pipelines. A major step forward in speech synthesis has been the incorporation of a diffusion model into TTS for the purpose of pitch prediction. The

diffusion model is a powerful tool for improving the expressiveness and realism of synthetic speech since it can produce a wide range of realistic pitch patterns and cope with the hierarchical and multi scale aspects of prosody.

The diffusion model is a powerful tool for text to speech systems, allowing them to generate natural sounding speech while still faithfully representing the input texts content and capturing the subtle prosodic fluctuations. Synthesized speech becomes more expressive, engaging, and resembling human speech patterns as a consequence. A diffusion model is a useful tool for many TTS applications, including pitch prediction, since it delivers practical advantages during model training and deployment. With the area of TTS constantly developing, the diffusion model's ability to improve pitch prediction and overall speech quality will undoubtedly play a larger and larger role. This will ultimately lead to speech synthesis systems that are more sophisticated and flexible.

$$q(X_t|X_{t-1}) := N(X_t|\sqrt{1 - \beta_t}X_{t-1}, \beta_t I) \quad (1)$$

During the forward operation of the diffusion process, the data is incrementally transformed into a totally noisy version by adding Gaussian noise. This step is fundamental because it lays the groundwork for the subsequent step that generates new data or restores the existing data. The forward process aims to gradually introduce noise into the data across several time steps until the data becomes completely noise like. Data noise onset rate is strongly influenced by the variance scheduling beta. When beta is modest, the noise is introduced more slowly, which means that the data keeps more of its original structure for a longer period of time throughout the forward process. Because it affects the reverse denoising process complexity and the quality of the produced data, the variance schedule choice may significantly affect the diffusion models' performance. When training a diffusion model to reverse a process, the forward process is an essential component.

At each stage of the reversal process, the model is trained by being exposed to different noisy versions of the data and its job is to forecast the original, clean data. Data creation makes advantage of the model's capacity to progressively reduce noise and reproduce the original data, which it acquires during training. This backwards process is more of a learnt procedure than a basic inversion; the model, which is often parameterized using deep neural networks, approximates the backward dynamics of the noise addition. A typical strategy is to develop the model with the objective of minimizing a loss function. This loss function measures the discrepancy between the predictions made by the model and the real data on iteration. The most often used approach is mean square error (MSE) loss algorithm. In this technique, the model makes an effort to accurately forecast the clean data from its noisy counterpart. As the model improves with practice, it may eventually produce high quality data by reversing the process and beginning with random noise.

By efficiently sampling from the learnt distribution of pitch patterns, this diffusion framework allows for the creation of various and natural pitch contours in the context of pitch prediction for TTS. This model is able to mimic human speech in its natural variety by using the reverse diffusion process on a random noise vector as a starting point. Applications like TTS, its objective is to produce speech that is as authentic and emotive as feasible, benefit greatly from the diffusion model's capacity to record and repeat these fluctuations. Important preliminary work in a diffusion model occurs during the forward phase, when the initial data is progressively turned into noise following a well-chosen variance schedule. After this noisy transformation, the model learns to reassemble the original data from its

noisy variants, which is the reverse process. Once this reversal process is mastered, the diffusion model transforms into a potent generative tool that can generate various, high-quality outputs, and such TTS friendly pitch contours.

$$q(X_{1:T} | X_0 := \prod_{t=1}^T q(X_t | X_{t-1}) \quad (2)$$

In which T is the total number of steps till the diffusion process ends. For sufficiently high values of T , we may treat x_T as random noise with a normal distribution $N(o, I)$.

$$p\theta(X_t - 1 | X_t) := N(X_t - 1; \mu\theta(X_t, t), \sigma_t^2 I) \quad (3)$$

The absence of an analytically straightforward method to derive the forward process that corresponds to the reverse process at each time step is one of the main issues in the context of diffusion models, especially those employed in generative tasks like pitch prediction for TTS. The model has to learn to approximate the forward process using a parameterized method because it is not immediately accessible. Here, the neural network represented here by the theta parameters comes into action. Every time step in the forward process is defined by a set of parameters that the neural network is taught to anticipate. Given the complicated and high dimensional nature of the data, the neural networks job is to provide a strong and versatile approximation of the forward process.

The neural network uses the noisy data that is available now to try to guess what the parameters would have been for the data that was free of noise at the prior time step. Through effective learning of the dynamics of the diffusion process, the network is able to replicate the progressive noise addition process even in the absence of information about the exact forward process. In order to optimize the parameters θ and train the neural network, the model usually minimizes a loss function that is constructed from the variational bound of the negative log likelihood. An effective method for approximating complicated distributions in probabilistic modelling; variational inference is the foundation of this approach. Minimizing the variational constraint guarantees that the learnt distribution of the model closely matches the genuine data distribution, even when the forward process cannot be directly formulated. In this case, the variational bound shines because it gives the neural network a manageable goal to optimize. In such a situation, the negative log likelihood quantifies the disparity among the actual records circulation and the model generated distribution.

The model learns to provide data that closely resembles the actual distribution by reducing this number, which enhances the output quality. The variational limit, more precisely, is a maximum allowable value for estimating the negative log probability of the data used in the model. The capacity of models to accurately depict the real distribution of data improves as the constraint decreases. The goal of the models training is to lower this limit so that it can better forecast the reverse process at each time step. When applied to tasks such as pitch prediction in TTS, this improves the model's ability to provide varied, high-quality outputs.

During training, the technique of reverse propagation is utilized to modify the model by calculating and using the slopes of a loss function in relation to the theta, which are the variables used by the neural network. To achieve a well-trained diffusion model that can generate realistic data and properly approximate the forward process, this iterative procedure keeps going until the model finds a set of parameters that minimize the variational constraint. The need of minimizing the variational bound during training is crucial. This mechanism is essential for the diffusion model as a whole, since it enables the neural network to develop a good approximation of the forward process. The model's

ability to provide accurate results depends on this approximation; without it, it would be unable to comprehend the addition of noise to the data throughout the forward phase.

The variational bound of the negative log likelihood is minimised during neural network training; this function stands in for the actual distribution of data. By going through this training process, the model may learn to provide varied and high-quality outputs by properly predicting the reverse process at each time step. Producing genuine and expressive speech relies on capturing the richness and unpredictability of the input, which is why this method is vital to diffusion models performance in tasks like pitch prediction for TTS.

$$L(\theta) := E_{t, X_0, \epsilon} [\|\epsilon - \epsilon\theta(\sqrt{a_t}X_0 + \sqrt{1-a_t}\epsilon_t)\|_2] \quad (4)$$

Where sample $\epsilon \sim N(0,1)$ and $\epsilon(\cdot)$ are the results that the neural network produces. According to our model, the denier estimates θ for which the criteria C_t and S are used. Since this is a model of the denoising process, the following equation describes it:

$$X_{t-1} = \frac{1}{\sqrt{a_t}} \left(X_t - \frac{1-a_t}{\sqrt{1-a_t}} \epsilon\theta(X_t, C_t, s, t) \right) + \sigma_t Z \quad (5)$$

Additionally, its inputs are subjected to gradient cutting. Hence, this models predictor is a standalone module. For text to speech (TTS) synthesis to work, it must be able to capture and reproduce prosodic details like timing, rhythm, and pitch fluctuations in order to produce expressive and lifelike speech. Meaning, emotion, and the speaker's identity are all profoundly impacted by prosody, which includes the rhythm, melody, and stress patterns of speech. Nevertheless, because to its hierarchical structure and the interaction between global consistency and local variability, effectively modelling prosody is a challenging endeavour.

Forecasting of the Prosody Model

The goal of prosody prediction is to generate speech that sounds natural by figuring out how to stress and intonate a series of phonemes. Speaking with the right prosody is essential for expressing one's feelings and ideas via language; speaking with the wrong prosody may make one seem robotic or bewildering. Take into consideration that the tone employed in the phrase You are going to the party might alter its meaning. The synthesis process in TTS systems is guided by prosodic contours generated by prosody prediction models, which examine language factors including sentence structure and word stress.

Phonemic context, sentence location, and the speaker's intended emphasis or moods are only a few of the interacting aspects that make prosody prediction fundamentally problematic. In their early days, TTS systems often relied on rule-based approaches to prosody prediction, using predetermined patterns determined by the phrase structure. Unfortunately, the subtleties of genuine speech were beyond the capabilities of these systems. Deep learning techniques are used by modern TTS systems to improve the accuracy of prosody prediction. These models may learn intricate patterns of prosodic variation since they are trained on massive voice data corpora. The models may improve the overall quality of the TTS output by including variables like pitch, duration, and loudness, which allow for more expressive and genuine speech generation.

RESULTS AND DISCUSSION

Finally, the researcher concluded that the diffusion model excels at handling the complex and diverse prosody of speech, making it an ideal choice for pitch prediction in TTS. Pitch changes in human speech may be attributed to a variety of factors, including but not limited to variances in emphasis, speaker identity, emotional state, and language context. Because of their inability to account for such variety, conventional pitch prediction methods often produce artificially produced speech that comes out as robotic or artificially produced. On the other hand, the diffusion model may generate more varied and realistic sounding pitch patterns due to its innate ability to generate various materials. In the first step, the encoded text and speaker representation are used as preconditions. Then, the pitch predictions are refined using the diffusion models reverse process. The final product is full of prosodic diversity since the pitch pattern is fine-tuned and complicated at each stage of the reversal process. Prosody, which includes the tempo, accent, and intonation of speech, is notoriously difficult to capture and reproduce in text to speech (TTS) synthesis, especially in zero shot situations when the model produces speech for invisible speakers. Synthesizing speech with naturalness and expressiveness relies heavily on prosody, which is equally important for communicating meaning and emotion in spoken language. This is an area where current TTS models often fail, especially when it comes to dealing with the complex prosodic changes across various language levels and keeping speakers consistent.

The diffusion model should be compared against traditional TTS models like Tacotron 2 + WaveGlow, FastSpeech 2 + HiFi-GAN, Grad-TTS (Diffusion-based), Your Proposed Model, Tabulated results for metrics like MOS, MCD, and WER would be insightful.

Table 1 Proposed model compares with existing

Model	MOS	MCD	WER	PESQ
Tacotron 2	3.85	4.92	6.30%	3.1
FastSpeech 2	3.9	4.8	5.90%	3.15
Grad-TTS(Baseline Diffusion)	4.1	4.5	5.10%	3.3
Proposed Model	4.35	4.2	4.50%	3.45

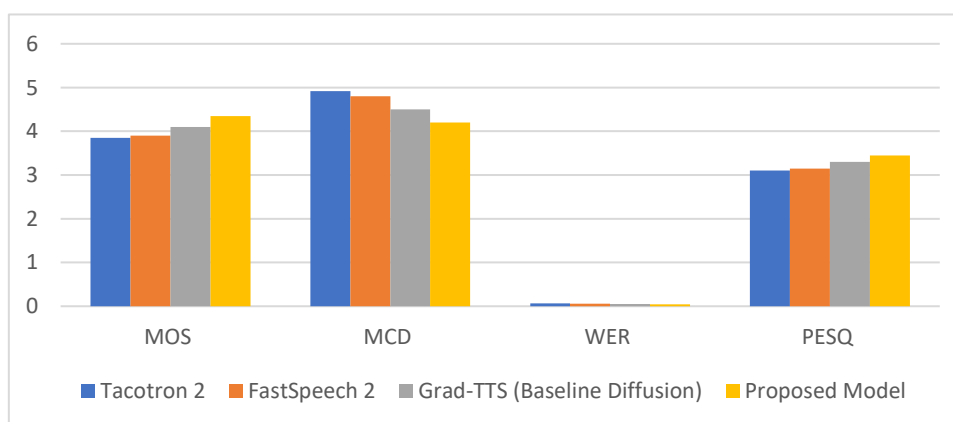


Figure 2. Proposed model result analysis.

CONCLUSION

The diffusion model is a better option for pitch prediction in TTS systems due to its remarkable ability to capture the complex prosodic aspects of speech. The model overcomes the drawbacks of conventional methods by generating legitimate, expressive, and distinct speech by iteratively improving pitch patterns through its reverse process. It works especially well in difficult situations like zero-shot synthesis, where prosody is essential to obtaining emotional resonance and naturalness. An important development in the realm of TTS synthesis is the diffusion model's capacity to manage intricate prosodic fluctuations while preserving speaker consistency.

REFERENCES

- [1] Indumathi A, Chandra E. Survey on speech synthesis. *Signal Process Int J*. 2012;6(5):140.
- [2] Jayaraman R, Vasanthi G, Ramaratnam MS. A study on investors' behavior towards equity and mutual funds. *Glob J Commer Perspect*. 2014;3(4):132–6.
- [3] John S, Chattopadhyay P. Factors impacting leadership effectiveness: A literature review. *Arab J Bus Manag Rev*. 2015; ISSN: 2223-5833.
- [4] *Journal of AI and Data Mining*. 2020;8(4):491–514.
- [5] Juvela L, Bollepalli B, Yamagishi J, Alku P. Waveform generation for text-to-speech synthesis using pitch-synchronous multiscale generative adversarial networks. In: *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing*; 2019 May; p. 6915–9.
- [6] Kaneko T, Kameoka H, Hojo N, Ijima Y, Hiramatsu K, Kashino K. Generative adversarial network-based postfilter for statistical parametric speech synthesis. In: *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing*; 2017 Jun; p. 4910–4.
- [7] Li N, et al. Neural speech synthesis with transformer network. In: *Proceedings of the AAAI Conference on Artificial Intelligence*; 2019;33(1).
- [8] Li Y, Qin D, Zhang J. Speech synthesis method based on Tacotron2. In: *Proceedings of the 13th International Conference on Advanced Computational Intelligence*; 2021; p. 94–9.
- [9] Lin S, Su W, Meng L, Xie F, Li X, Lu L. Nana-HDR: A non-attentive non-autoregressive hybrid model for TTS. *arXiv Preprint arXiv:2109.13673*. 2021.
- [10] Ling ZH, Deng L, Yu D. Modeling spectral envelopes using restricted Boltzmann machines for statistical parametric speech synthesis. In: *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing*; 2013; p. 7825–9.
- [11] Ling ZH, Deng L, Yu D. Modeling spectral envelopes using restricted Boltzmann machines and deep belief networks for statistical parametric speech synthesis. *IEEE Trans Audio Speech Lang Process*. 2013;21(10):2129–39.
- [12] Mohamed AR, Dahl GE, Hinton G. Acoustic modeling using deep belief networks. *IEEE Trans Audio Speech Lang Process*. 2011;20(1):14–22.
- [13] Moradi M, Hamidzadeh J. Ensemble-based top-k recommender system considering incomplete data. *J AI Data Min*. 2019;7(3):393–402.
- [14] Ravi DJ, Patilkulkarni S. Text-to-speech synthesis system for Kannada language. *Int J Adv Res Comput Sci*. 2021;2(1):298–304.

- [15] Rebai I, Ben Ayed Y. Arabic text-to-speech synthesis based on neural networks for MFCC estimation. In: Proceedings of the World Congress on Computer and Information Technology; 2013 Jun; p. 1–5.
- [16] Ren Y, et al. FastSpeech: Fast, robust and controllable text-to-speech. Adv Neural Inf Process Syst. 2019;32.
- [17] Ren Y, Tan X, Qin T, Zhao S, Zhao Z, Liu TY. Almost unsupervised text-to-speech and automatic speech recognition. In: Proceedings of the International Conference on Machine Learning; 2019; p. 5410–9.
- [18] Saito Y, Takamichi S, Saruwatari H. Statistical parametric speech synthesis incorporating generative adversarial networks. IEEE/ACM Trans Audio Speech Lang Process. 2017;26(1):84–96.
- [19] Tan X, Ren Y, He J, Zhou Z, Qin T. NaturalSpeech: End-to-end text-to-speech synthesis with human-level quality. arXiv Preprint arXiv:2205.04425. 2022.
- [20] Thu CST, Zin T. Implementation of text-to-speech conversion. Int J Eng Res Technol (IJERT). 2020;3(3).
- [21] Toma SA, Tarsa GI, Oancea E, Munteanu DP, Totir F, Anton L. A TD-PSOLA based method for speech synthesis and compression. In: Proceedings of the 8th International Conference on Communications; 2010; p. 241–50.
- [22] Van den Oord A, Dieleman S, Zen H, Simonyan K, Vinyals O, Graves A, et al. WaveNet: A generative model for raw audio. arXiv Preprint arXiv:1609.03499. 2016.
- [23] Wagner D, Bayerl SP, Cordourier Maruri HA, Bocklet T. Generative models for improved naturalness, intelligibility, and voicing of whispered speech. In: Proceedings of the IEEE Spoken Language Technology Workshop; 2022; p. 943–8.
- [24] Wang G. Deep text-to-speech system with seq2seq model. arXiv Preprint arXiv:1903.05955. 2019.
- [25] Wang X, Takaki S, Yamagishi J. An RNN-based quantized FO model with multi-tier feedback links for text-to-speech synthesis. In: Proceedings of INTERSPEECH; 2017; p. 1059–63.
- [26] Watts O, Stan A, Clark RA, Mamiya Y, Giurgiu M, Yamagishi J, King S. Unsupervised and lightly-supervised learning for rapid construction of TTS systems in multiple languages from found data: Evaluation and analysis. In: Proceedings of the Eighth ISCA Workshop on Speech Synthesis. 2013.
- [27] Gopi A, Sajini T, Bhadrar VK. Implementation of Malayalam text to speech using concatenative based TTS for Android platform. In: Proceedings of the International Conference on Control Communication.
- [28] Gopi A., Sajini T., Bhadrar V.K. "Implementation of Malayalam text to speech using concatenative based TTS for Android platform." In: Proceedings of the International Conference on Control Communication. 2015.
- [29] Ravi, D. J., Patilkulkarni, S. "Text-to-speech synthesis system for Kannada language." International Journal of Advanced Research in Computer Science. 2021;2(1):298-304.
- [30] Madhusree R. "Design and implementation of a TTS system for Indian languages." International Journal of Computer Applications. 2013;79(6):1-6.
- [31] Sreenivas, R., Rao, B. "Concatenative speech synthesis for Indian languages." International Journal of Speech Technology. 2017; 20(2):191–199.