

Unlocking the Potential: an evaluation of Text-to-Speech Models for the Bahnar Language

Giang Dinh Lu
University of Social Sciences
and Humanities, Vietnam
Vietnam National University
Ho Chi Minh City,
Ho Chi Minh City, Vietnam
giangdl@hcmussh.edu.vn

Hai Vu Hoang
Ho Chi Minh City University
of Technology (HCMUT),
Vietnam
Vietnam National University
Ho Chi Minh City,
Ho Chi Minh City, Vietnam
hai.vu.tharios19@hcmut.edu.vn

Tho Quan Thanh
Ho Chi Minh City University
of Technology (HCMUT),
Vietnam
Vietnam National University
Ho Chi Minh City,
Ho Chi Minh City, Vietnam
qttho@hcmut.edu.vn

Quy Nguyen Tran
University of Social Sciences
and Humanities, Vietnam
Vietnam National University
Ho Chi Minh City,
Ho Chi Minh City, Vietnam
tranquynghien@hcmussh.edu.vn

Duc Nguyen Quang
Ho Chi Minh City University
of Technology (HCMUT),
Vietnam
Vietnam National University
Ho Chi Minh City,
Ho Chi Minh City, Vietnam
nqduc@hcmut.edu.vn

Abstract

The paper aims at evaluating the effectiveness of an AI based mobile application of text- to-speech models for Bahnar language. In this application, a sequential combination of two models was implemented, starting with the application of the Grad-TTS model and subsequently followed by the Hifi-GAN model. Grad-TTS was employed to ensure a highly correct pronunciation of Bahnar words without being constrained by the dataset. The strengths of Hifi-GAN, in other hands, have been fine-tuned for the Bahnaric language to enhance the quality of synthesized audio, in order to produce a native-like Bahnar voice and accent. Those artificially generated sounds from our model achieved a high level of naturalness.

Keywords

Bahnar language, speech synthesis, text-to-speech conversion, Mean Opinion Score (MOS), modified rhyme test

I. INTRODUCTION

Languages are a crucial medium for communication and information transmission, especially in a multi-ethnic country like Vietnam. The diversity of languages here lead to a growing need for an effective machine translation systems. Machine translation therefor becomes more and more an important tool to assist people in multilingual communication and diverse settings of mass media among various ethnic communities.

The conversion of text into spoken form becomes progressively necessary and is seen to be an economical solution while Bahnar language is widely used in mass media (TV, radio, school, cultural events, training workshops...) as well as in different language domains, such as in commerce, business, education, and local activities. In Vietnam, the Bahnar diaspora consists of approximately 287,000 people (according to the 2019 Vietnam Population Census). They reside mainly in three provinces: Kon Tum, Gia Lai, and Binh Dinh, with different dialects maintaining mutual comprehensibility. The dialectal variation becomes challenging as producing accurate and natural Bahnar language audio converted from text needs to overcome the complexity in its syllabic structure, phonemes, and intonations.

Despite many research efforts and developments in the field of Text-to-Speech (TTS) in Vietnam, Bahnar TTS,

among most of ethnic languages, is entirely a new task. Conventional TTS models based on rules and statistics have considerable limitations in accuracy and in naturalness of synthetic audio. The advancement of deep learning models, especially Deep Neural Networks (DNN) and Recurrent Neural Networks (RNN), has created opportunities to enhance the efficiency of machine translation.

In this research, a fundamental research question was raised: Can a deep neural network-based text-to-speech model for the Bahnar language improve efficiency compared to the traditional methods? The current research focuses on evaluating and comparing the effectiveness of modern deep learning-based machine translation models in creating audio for Bahnar language. The aim of the method is to achieve better, more natural sound generation for various real-world applications.

In Vietnam, since 2010, researchers have been applying Text-to-Speech synthesis technology to the Vietnamese language. Võ Quang Diệu Hà and colleagues developed a speech synthesis system based on concatenative techniques [1]. They used syllables as the basic unit to create speech in their TTS system. A primary limitation of this method is the complex selection of speech database units by the system. Furthermore, the final product often doesn't express intonation and rhythm well.

Nguyễn Trang constructed a Vietnamese TTS system, developed on the Mary TTS platform, using a statistics-based speech synthesis method [2] with an architecture based on the Hidden Markov model [3]. The research results indicated a specific need for further enhancement in processing tonal aspects to generate better automated speech for Vietnamese language.

In 2019, Lâm Phùng Việt and his colleagues utilized a Deep Learning model to improve a TTS system [4]. The system is based on the Tacotron 2 model and the WaveGlow neural vocoder. Comparing it to the previous SPSS approach, the results show significantly higher performance.

In 2020, Chenfeng Miao utilized the Tacotron 2 model and the Hifi-gan vocoder to synthesize speech. The research results showed that the synthesized sound achieved up to 89.3% similarity with the natural speaker's voice [5].

However, there still exists a significant challenge when processing long sentences, especially in expressing intonation and accurately pronouncing borrowed words.

Vadim Popov suggested that the Grad-TTS model, when applied in TTS systems, can address the limitations of previous models, such as the requirements for extensive data, unnatural sound, lack of intonation, and difficulty handling long passages. Currently, the Grad-TTS model has shown positive results in English. As the Bahnar language lacks the intonation found in English and has a simpler syllabic structure, there is hope that it will yield positive results in synthesizing Bahnar speech.

Classical methods used to build Text-to-Speech systems include articulatory synthesis. However, there still exists a significant challenge when processing long sentences, especially in expressing intonation and accurately pronouncing borrowed words.

Vadim Popov suggested that the Grad-TTS model, when applied in TTS systems could indicate potential limitations of the previous models, such as requirements for extensive data, unnatural sound production, lack of intonation, and difficulty handling long passages. Currently, the Grad-TTS model has shown positive results in English. As the Bahnar language lacks the intonation found in English and has a simpler syllabic structure, there is hope that it will yield positive results in synthesizing Bahnar speech.

Classical methods were known to build TTS systems include articulatory synthesis [6], formants synthesis [7], concatenative synthesis [8] and statistical parametric speech synthesis [9]. Applying advanced techniques and achieving high performance in languages with limited data, such as Bahnar, is a challenge. We'll describe the process of selecting and building a suitable AI model, as well as constructing a suitable phoneme set for the Bahnar language. Therefore, this report also serves as a reference document and is one of the first research works in developing an artificial speech synthesis system for Bahnar.

The main goal of TTS systems is to convert arbitrary text into spoken language in the form of audio. Text processing and speech synthesis are the two main components of the text-to-speech system. The aim of text processing is to analyze the input text and generate phoneme sequences. These phonemes are marked by the speech synthesis component or are synthesized from parameters collected from a sufficiently large audio data source. To produce natural speech synthesis, the text processing is structured in an appropriate sequence, matching the phonemes to the arbitrary input text.

Currently, several models have been applied for text-to-speech conversion. However, selecting an appropriate model for an isolated tone language like Vietnamese requires specific attention. This becomes particularly crucial for Bahnar, a language lacking tonality. In terms of grammatical structure and vocabulary, Vietnamese and Bahnar share many similarities.

II. METHOD

One of the prominent methods applying AI to TTS is Tacotron 2 [10], utilizing a combined structure of Recurrent Neural Networks (RNN) and Convolutional Neural Networks (CNN). Tacotron 2 has become more and more popular in this direction. However, Tacotron 2 were evaluated by many as not meeting the requirements for natural speech synthesis. The

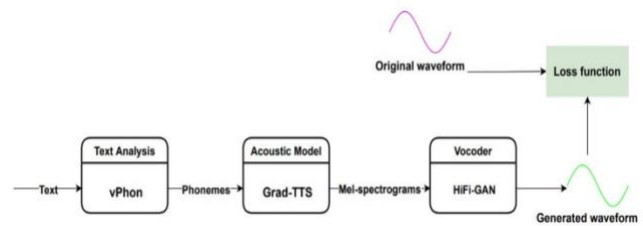
input text (INPUT) after processing by the analysis module results in a corresponding phoneme sequence (PROCESSED). This character string also serves as input for the training and utilization of the Grad-TTS model.

INPUT: Inh kăt 'ba kópung adoí.

PROCESSED: I-ɲ-k-ă-t- 6-a-k- ə-p- i- ɲ-a-d- ə-j.

(English: I'm harvesting rice in the field.)

The input can be a sentence or a paragraph. The computer will analyze each word in the sentence into phonemes. Based on the sound database for phonemes, which includes consonants, vowels, and semi-vowels, the computer will synthesize these phonemes to reproduce the sound. The sound produced by the model is not simply a replay of previously recorded sound but a synthesis process from phonemes in the



Bahnar language.

Sound quality is maintained by recording speech in a naturally noisy environment. The recorded audio data is then converted to digital data using software like Wavsurfer. The current system uses 16-bit resolution data along with a sampling rate of 22.05 KHz. (A sampling rate of around 20 KHz is considered sufficient to maintain sound quality). Recording of each phoneme is done at a 22.05 kHz sampling rate. The phoneme sound database includes all consonants and vowels.

We utilize the conversion method of vPhon [12] The rule of phoneme analysis for the Vietnamese and Bahnar transcription. The computer learns from improved phonemes because they represent sound more closely. Moreover, vPhon is specifically designed for Vietnamese. For the most part, the phonetic learning structure corresponds between Vietnamese and Bahnar. Bahnar has the advantage of resource-saving without intonation.

The second module in the process is the acoustic model based on Grad-TTS. The input to this model is the pseudo-phoneme set, and through the training process, it generates a mel-spectrogram. A mel-spectrogram is a spectral representation of sound waves, encompassing frequency and time. It can reproduce detailed information about frequency bands predominating at each moment in the sound wave. From the mel-spectrogram, you can extract the original sound wave through the inverse problem.

The final step is the vocoder processing, where you use the HiFi-GAN network to generate Bahnar sound. What's unique about a vocoder is that it directly calculates on the mel-spectrogram, thus not depending on the input language. Therefore, HiFi-GAN can be directly applied to Bahnar without retraining. The results show that HiFi-GAN still achieves high-quality sound compatible with the Bahnar language, making the synthesis of Bahnar speech efficient and high-quality.

III. EXPERIMENTS

A. Dataset interpretation

Dataset information

A dataset of 10,000 speech segments from a single speaker, each 3-7 seconds long with corresponding phonemes, was collected and labeled with unique IDs that correspond to the names of the corresponding .wav files. The contributor's voice is clear, loud, expressive, and low-noise, with minimal variation in delivery or tone.

Audio segmentation

In the beginning, audiobooks that last approximately 20 hours are first selected. Then, they are exported to .wav format using the pydub library. On later version of pydub, we use the split on silence to segment the audio based on pauses. It is worth noting the two most important parameters of evaluating the freshness of data, which are *minimum_silence_length* (in milliseconds), which is the minimum duration of silence used to split the audio, and *silence_thresh* (in dBFS), which is the threshold below which sound is considered silence (the default is -16 dBFS). Choosing incorrect values for these parameters can lead to noisy data. Therefore, the best value can be determined by examining the waveform of an audio clip to identify the speaker's segmentation pattern. After splitting the audio, we collect all speech segments that are between 3 - 7 seconds long and use the cognitive service from Azure to generate text transcripts for each speech segment.

Preprocessing

Outliers are removed from the dataset by plotting a scatter graph of the duration of each audio record. Linear regression is then applied using two variables: duration and the number of words in each file, since the number of words in a phrase is proportional to the audio file's duration. Outliers on the graph are eliminated.

Next, the text is normalized by converting numbers, ordinals, and currency units into complete words (UTF-8). All punctuation is then removed from the text to enable the model to learn segmentation independently.

Finally, a phonemic analysis is generated for each audio record. This is the second attribute of the dataset used for model training. The transformation rules of vPhon [18] are applied to analyze phonemes for each audio record. The machine learns better from phonemes than from the original text because phonemes are representations closer to sound.

B. Bahnar synthesis model

The model is developed on the PyTorch platform. Overall, the data is throughput for over 2500 iterations on a server equipped with an NVIDIA RTX 3090 GPU. The accompanied audio recordings, with a total duration of approximately 57 hours, were divided into smaller segments to better serve the training and evaluation process. Before being fed into the TTS system, all audio files underwent filtering and preprocessing process to ensure no noise residues nor distortions. In addition, all input text was analyzed and converted into phoneme forms to generate suitable data for the model.

DCA-Tacotron 2 model is the forerunner on publicly available data (InfoRe dataset), using the same parameters and configurations as in the original paper. Subsequently, we evaluated the effectiveness of our proposed data enrichment

method and framework on both the InfoRe dataset and our ViSpeech dataset, comparing it with the baseline model. Conducting an exploratory experiment on the Bahnar language, a low-resource language, was intended to demonstrate the utility of our proposed approach. For each dataset, 90% of the samples are used for training and the remaining 10% are reserved for validation. Previously, we selected and set aside 40 sentences for testing.

IV. RESULTS AND DISCUSSION

Once the Bahnar audio synthesis process was completed, the natural voices of Bahnar collaborators from the Gia Lai region were directly compared with the voices generated by the model. To evaluate the reliability of the artificial Bahnar sound, the Comparative Mean Opinion Score (CMOS) test was utilized. With such test, a total of 60 pairs of sentences with corresponding audio were randomly selected from the dataset. The input text was then fed into the model to generate the synthetic Bahnar voice. This artificial Bahnar voice was then directly compared to the natural Bahnar voice..

To evaluate the quality of the artificial Bahnar voice, 30 Bahnar informants from three areas: Gia Lai, Kon Tum, and Binh Dinh, were invited to participate in the assessment. For each dialect (Gia Lai, Kon Tum, and Binh Dinh), three sets of assessment samples were created. Each collaborator listened to both versions of the voice (from a real person and the model-generated one) and then assessed the similarity level on a scale ranging from -3, meaning the model-generated voice is much worse than the real voice, to +3, meaning vice-versa. The obtained scores were then used to evaluate the quality of the artificial Bahnar voice.

Finally, the average scores were calculated and displayed in Table 1. The results indicate that the artificial intelligence model is only slightly inferior to the human real voice. Based on this outcome, it can be concluded that the selected model has achieved high performance and quality for the proposed artificial voice generation system.

Table 1: CMOS results comparing real voice and neural network-generated voice."

System	CMOS
Ground truth	0.000
OURS	-0.3276

Furthermore, we have developed two additional versions, one based on a female reading and another based on both male and female voices. Both models have yielded promising results and are intelligible to the Bahnar people. This implies that we have expanded the model's application to accommodate both male and female voices, enabling the creation of high-quality sound for various purposes and user groups within the Bahnar community.

V. CONCLUSIONS

In this study, we established the initial foundation for developing a natural speech synthesis technology for the

Bahnar language in Vietnam. Despite the language's unique characteristics and limited dataset, our approach yielded promising and effective results. The speech generated by our system closely mimicked human speech, exhibiting flexibility in handling input text. The speech quality from our system was evaluated as clear and easily understandable, able to simulate both male and female voices as required.

Evaluation results demonstrated that the system is capable of producing high-quality sound, and the voice conversion model trained on Grad-TTS source data combined with the HiFi-GAN vocoder showed superior performance. Future work may involve enhancing the quality of noise-disturbed sound.

ACKNOWLEDGMENT

This research is funded by Ministry of Science and Technology (MOST) within the framework of the Program "Supporting research, development and technology application of Industry 4.0" KC-4.0/19-25 – Project "Development of a Vietnamese- Bahnaric machine translation and Bahnaric text-to-speech system (all dialects)" - KC-4.0-29/19-25.

REFERENCES

- [1] V. Q. D. Ha, N. M. Tuan, C. X. Nam, P. M. Nhut, and V. H. Quan, "Vos: The corpus-based vietnamese text-to-speech system," 2010.
- [2] H. Zen, K. Tokuda, and A. W. Black, "Statistical parametric speech synthesis," *Speech communication*, vol. 51, no. 11, pp. 1039-1064, 2009.
- [3] L. Rabiner and B. Juang, "An introduction to hidden Markov models," *Eee assp magazine*, vol. 3, no. 1, pp. 4-16, 1986.
- [4] V. L. Phung, H. K. Phan, A. T. Dinh, K. D. Trieu, and Q. B. Nguyen, "Development of Zalo Vietnamese Text-to-Speech for VLSP 2019."
- [5] Tung Tran *et al.*, "Naturalness improvement of Vietnamese Text-to-Speech System using Diffusion Probabilistic modelling and Unsupervised Data Enrichment," in *The First International Conference on Intelligence of Things (ICIT 2022)*, 2022.
- [6] A. B. L. Larsen, S. K. Sønderby, H. Larochelle, and O. Winther, "Autoencoding beyond pixels using a learned similarity metric," in *International conference on machine learning*, 2016, pp. 1558-1566: PMLR.
- [7] I. Goodfellow *et al.*, "Generative adversarial nets," in *Advances in neural information processing systems*, vol. 27, M. I. Jordan, Y. LeCun, and S. A. Solla, Eds., 2014.
- [8] X. Mao, Q. Li, H. Xie, R. Y. Lau, Z. Wang, and S. Paul Smolley, "Least squares generative adversarial networks," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2794-2802.
- [9] K. Kumar *et al.*, "Melgan: Generative adversarial networks for conditional waveform synthesis," in *Advances in neural information processing systems*, vol. 32, 2019.
- [10] J. Shen *et al.*, "Natural tts synthesis by conditioning wavenet on mel spectrogram predictions," in *2018 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, 2018, pp. 4779-4783: IEEE.
- [11] V. Popov, I. Vovk, V. Gogoryan, T. Sadekova, and M. Kudinov, "Grad-tts: A diffusion probabilistic model for text-to-speech," in *International Conference on Machine Learning*, 2021, pp. 8599-8608: PMLR.
- [12] J. Kirby, "VPhon: A Vietnamese Phonetizer," ed, 2008.