

Voice conversion for natural-Sounding speech generation on low-Resource languages: A case study of bahnaric

Dang Tran Dat^{1,2}, Tang Quoc Thai^{1,2}, Duc Q. Nguyen^{1,2}, Vo Duy Hung^{1,2}, Quan Thanh Tho^{1,2,*}



Use your smartphone to scan this QR code and download this article

ABSTRACT

Bahnar is an ethnic minority group in Vietnam, prioritized by the government for the preservation of their cultural heritage, traditions, and language. In the current era of AI technology, there is substantial potential in synthesizing Bahnar voices to support these preservation endeavors. While voice conversion technology has made strides in enhancing the quality and naturalness of synthesized speech, its focus has predominantly been on widely spoken languages. Consequently, low-resource languages like the Bahnaric language family encounter numerous disadvantages in voice synthesis. This study addresses the formidable challenge of synthesizing natural-sounding speech in low-resource languages by exploring the application of voice conversion techniques to the Bahnaric language. We introduce the BN-TTS-VC system, a pioneering approach that integrates a text-to-speech system based on Grad-TTS with voice conversion techniques derived from StarGANv2-VC, both tailored specifically for the nuances of the Bahnaric language. Grad-TTS allows the system to articulate Bahnaric words without vocabulary limitations, while StarGANv2-VC enhances the naturalness of synthesized speech, particularly in the context of low-resource languages like Bahnaric. Moreover, we introduce the Bahnaric-fine-tuned HiFi-GAN model to further enhance voice quality with native accents, ensuring a more authentic representation of Bahnaric speech. To assess the effectiveness of our approach, we conducted experiments based on human evaluations from volunteers. The preliminary results are promising, indicating the potential of our methodology in synthesizing natural-sounding Bahnaric speech. Through this research, we aim to make significant contributions to the ongoing efforts to preserve and promote the linguistic and cultural heritage of the Bahnar ethnic minority group. By leveraging the power of AI technology, we aspire to bridge the gap in speech synthesis for low-resource languages and facilitate the preservation of their invaluable cultural heritage.

Key words: Bahnaric speech synthesis, text-to-speech, natural-sounding voice conversion

¹Faculty of Computer Science and Engineering, Ho Chi Minh City University of Technology (HCMUT), Vietnam

²Vietnam National University Ho Chi Minh City (VNU-HCM), Vietnam

Correspondence

Quan Thanh Tho, Faculty of Computer Science and Engineering, Ho Chi Minh City University of Technology (HCMUT), Vietnam

Vietnam National University Ho Chi Minh City (VNU-HCM), Vietnam

Email: qttho@hcmut.edu.vn

History

- Received: 08-9-2023
- Accepted: 27-3-2024
- Published Online:

DOI :



Copyright

© VNUHCM Press. This is an open-access article distributed under the terms of the Creative Commons Attribution 4.0 International license.



1 INTRODUCTION

2 The Bahnar or Ba-Na (Vietnamese pronunciation:
3 [ba˧na˧]) represents a distinct ethnic minority within
4 the diverse tapestry of ethnic populations in Viet-
5 nam. Contemporary efforts spearheaded by the Viet-
6 namese government aim to enhance their integra-
7 tion through advancements in socio-cultural and sci-
8 entific literacy. A significant portion of this en-
9 deavor includes translating key documents into the
10 Bahnaric language by governmental and community
11 stakeholders. Concurrently, there is growing inter-
12 est among domestic research groups to devise auto-
13 matic translation systems for Vietnamese to Bahnaric
14 ethnolects. Notwithstanding these advancements, the
15 distinct characteristics of the Bahnar, given their sta-
16 tus as a smaller ethnic faction, result in hesitations
17 in engaging with the predominant Kinh (Vietnamese)
18 population. This occasionally impedes their complete
19 access to written information. Thus, conveying infor-
20 mation with native-like Bahnaric speech could signif-

icantly enhance accessibility for this community. 21
Modern TTS (text-to-speech) systems¹ can assist in 22
pronouncing words from text based on a trained 23
dataset. However, these systems require a substan- 24
tial amount of training data. For extremely low- 25
resource languages like Bahnaric, gathering a high- 26
quality training dataset becomes particularly arduous, 27
resulting in suboptimal pronunciation outputs. For 28
the small Bahnaric ethnic group, this also poses sig- 29
nificant challenges to communication. 30
Another solution is to develop voice conversion sys- 31
tems² that convert the voice quality to match that 32
of a genuine Bahnar individual. Due to the low- 33
resource nature of Bahnaric, we have proposed an ef- 34
fective approach that combines the Grad-TTS model³ 35
and the StarGANv2-VC model⁴. The use of the 36
Grad-TTS model enables the system to pronounce 37
an unlimited vocabulary from available texts. Mean- 38
while, the StarGANv2-VC model assists in gener- 39
ating a converted voice from an existing Bahnaric 40

Cite this article : Dat D T, Thai T Q, Nguyen D Q, Hung V D, Tho Q T. **Voice conversion for natural-Sounding speech generation on low-Resource languages: A case study of bahnaric.** *Sci. Tech. Dev. J. – Engineering and Technology* 2024; ():1-12.

voice. Particularly, the combination of Grad-TTS and StarGANv2-VC aids in refining and cleaning words and phonemes that Grad-TTS has not generated well, especially when trained from low-resource and low-quality sources like direct recordings of Bahnaric people's speech. In addition, we also introduce the HiFi-GAN-BN model, a variant of HiFi-GAN⁵ pre-trained by Bahnaric voice, to resemble the Bahnaric accents better when transforming the mel-spectrogram output of StarGANv2-VC into human-listenable waveform.

We have experimented with our system, known as BN-TTS-VC, using real-world data collected from the Bahnar community in the provinces of Gia Lai, Kon Tum, and Binh Dinh. When evaluated by human assessments, we have obtained favorable results.

The remainder of the paper is organized as follows. Section 2 describes previous works which are related to our study. Section 3 gives details of the Bahnaric phonological system. Section 4 describes the methodology to develop the BN-TTS-VC system. Section 5 presents the experiment results. Section 6 provides a discussion of the results obtained from our experiment. Section 7 presents conclusions and future work.

RELATED WORKS

Text-to-speech techniques

Text-to-speech synthesis is a task that involves converting written text into spoken words. The goal is to generate synthetic speech that sounds natural and resembles human speech as closely as possible. Classical methods used to construct text-to-speech systems include articulatory synthesis⁶, formant synthesis⁷, concatenative synthesis⁸, and statistical parametric speech synthesis⁹. These methods usually generate a voice with less of a natural or lack of emotion and the voice quality is low due to containing screeching and jerking sounds. Certain end-to-end models such as ClariNet¹⁰, FastSpeech 2s¹¹, and EATS¹² that create audio directly from text have been proposed based on simplification of text analysis modules and directly taking character strings or phonemes as input, also as to simplify acoustic properties with timbre spectra. The advantages of neural network-based speech synthesis over previous Text-to-speech systems include high voice quality in terms of intelligibility and naturalness as well as less reliance on the construction of input properties. Concerning Vietnamese text-to-speech systems, the Tacotron 2 acoustic model¹³ is considered a classical deep-learning method that is widely applied in these systems. The ZALO group developed a Text-to-speech system¹⁴ based on Tacotron

2¹³ and WaveGlow¹⁵ whose performance of their system is superior to the statistical parametric speech synthesis classical method.

Voice conversion techniques

Voice conversion (VC) is a technique for converting one speaker's voice identity into another while preserving linguistic content. Though most voice conversion methods that require parallel utterances achieve high-quality natural conversion results, it strongly limits the conditions to apply. Regarding non-parallel voice conversion methods, it can mainly be divided into three categories. *Auto-encoder approach*¹⁶⁻¹⁹ requires carefully designed constraints to remove speaker-dependent information, and the converted speech quality depends on how much linguistic information can be retrieved from the latent space. By contrast, *GAN-based approaches*, such as CycleGAN-VC3²⁰ use a discriminator that teaches the decoder to generate speech that sounds like the target speaker. Due to the lack of learning meaningful features from the real data in the discriminator, this approach often suffers from problems such as dissimilarity between converted and target speech, or distortions in voices of the generated speech. On the other hand, *TTS-based approaches* like Cotatron²¹, AttS2S-VC²², and VTN²³ extract aligned linguistic features from the input speech to give the converted speaker identity that is similar to the target speaker identity. However, the text labels for this approach are not often available at hand.

BAHNARIC HONOLOGICAL SYSTEM

To develop a speech synthesis system, it is essential to construct a phonological system for this particular language. Figure 1 illustrates an example of a Bahnaric language text. We can see that the language has its characteristics, and using the input parsing modules of other languages is impossible. Therefore, we analyze this language elaborately and build a set of pseudo-phonemes for the Bahnaric language, which is suitable input for the text-based speech generation model. The set of pseudo-phonemes is shown in Figure 2.

Each word in the input text will be compared to the corresponding phoneme sequence based on the above alphabet. An example is shown in Figure 3. From the text (INPUT) passed through the analyzer, the result is the corresponding phoneme sequence (PROCESSED). That sequence is also the input for training and using the TTS model.

adriêng nganh y teâ adriêng bet teêk weêk pôloêk phun bôgang bet sôhmeêch
 minh suaât kua tri giaê 01 trieâu ñoàng
 tôplih lôêm tôdrong tôme rong jaêng pran ñeh oei xa vinh kim
 trô jeân pôm minh sônaêm kung thu yoêk ñei khoang 60 trieâu ñoàng
 rim mô hình anu jôh pôjing thu yoêk tôpaê pônhoâm lö naê ma adriêng pôm

Figure 1: An example of text in Bahnaric language.

a b c d e f g h i j k l m n o p q r s t u v w x y z à á â ã ä å è é ê ì í î ñ ò ó ô õ ö ø ù ú û ý ã ĩ đ ũ ŕ ạ á ả à ă ằ ẳ ẵ ề ề ể ệ ì ï ọ ồ ờ ỗ ộ ớ ờ ỡ ợ ừ ừ ữ ự

a) *Monophonic*
 ia iă ie iě iô iö ua uă ue uě uê

b) *Diphthong vowels*

b l br by ch dj dr gl gr gy hl hm hn hñ hr hy jr kh kl kr ky ly ml mr ny my ñr ng ph pl pr py sr th tr ty

c) *Double consonants*
 hng ngl nhr

d) *Triple consonants*

Figure 2: A set of pseudo-phonemes for Bahnaric language.

INPUT: adriêng nganh y teâ
 PROCESSED: a-d-r-i-ê-ng ng-a-n-h y t-e-â

Figure 3: An example of an input text analyzer in Bahnaric language.

142 RESEARCH METHODOLOGY

143 Overview of the combined system of Text-to-speech and voice conversion for Bahnaric language

146 This system is constructed based on two main mod- 159
 147 ules including Text-to-speech and Voice Conversion, 160
 148 as illustrated in Figure 4. The first module gets the 161
 149 Bahnaric language text as input to generate a native 162
 150 voice with the content of the input text. There are
 151 two sub-models in this module, which are the vocoder
 152 and acoustic model. While the acoustic model gener-
 153 ates acoustic properties directly from input phonemes
 154 mentioned in Section 3, a vocoder transforms these
 155 features into sound waveforms. After that, the sound
 156 waveforms are passed to the Voice conversion module
 157 for generating the other types of voice of native based
 158 on the reference voice. This module is built from three

main component models for the purpose of extract- 159
 ing the characteristics of voice, converting the voice, 160
 and transforming the mel-spectrogram into a human- 161
 listenable waveform. 162

Grad-TTS system for Bahnaric speech synthesis

163 According to our research, there so far has been no 164
 165 reported work on building an artificial voice gener- 166
 166 ation system for the Bahnaric language. In this do- 167
 167 main, there is an existence of certain different char- 168
 168 acteristics between the Bahnaric and other popular 169
 169 languages. Therefore, applying techniques with high 170
 170 efficiency in those languages to Bahnaric is a highly 171
 171 complex problem. 172

173 One of the typical methods of applying AI to solve 174
 174 this problem is Tacotron 2¹³, which uses the archite- 175
 175 cture of recurrent neural network (RNN) and convo-

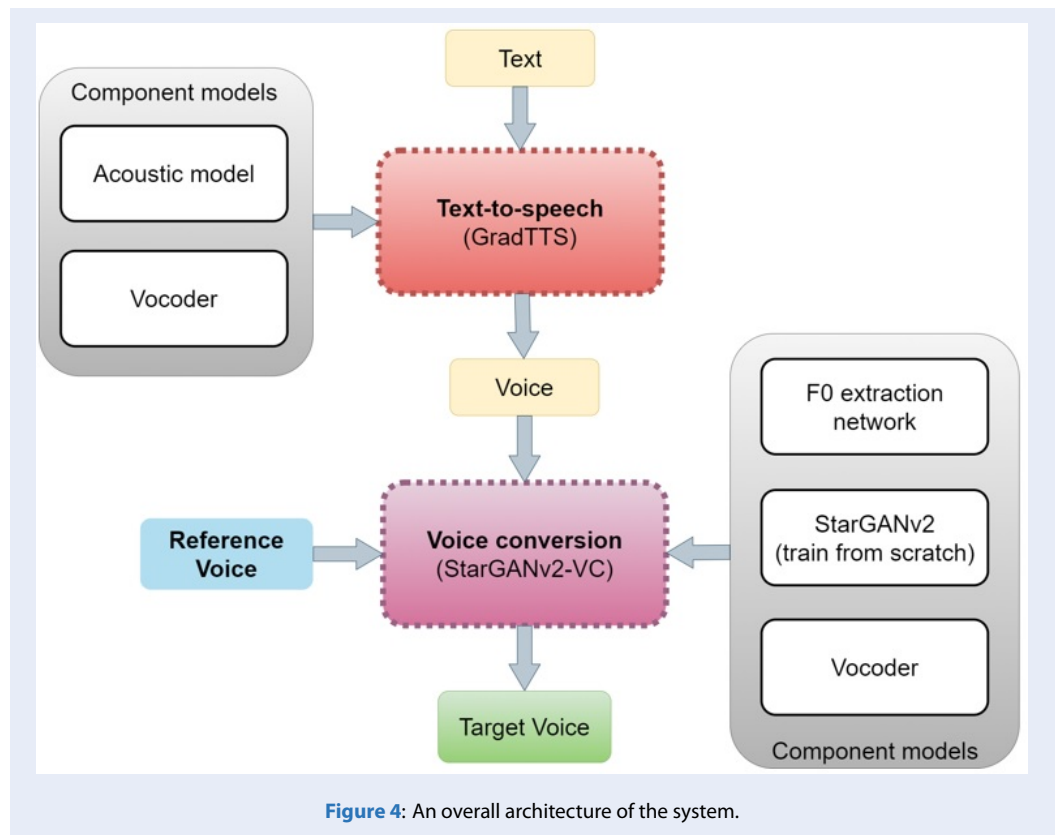


Figure 4: An overall architecture of the system.

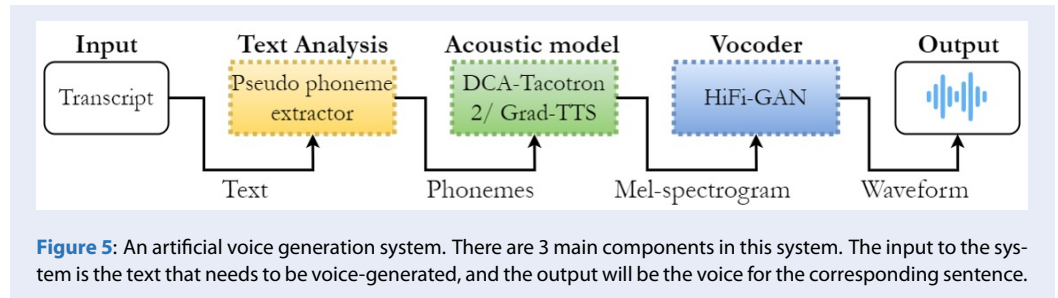


Figure 5: An artificial voice generation system. There are 3 main components in this system. The input to the system is the text that needs to be voice-generated, and the output will be the voice for the corresponding sentence.

176 lutional neural network (CNN). Tacotron 2 has been
 177 and is being used commonly in this field. However,
 178 this model has not yet met the requirements for the
 179 naturalness of the generated voice, especially for lan-
 180 guages such as Vietnamese and Bahnaric, because
 181 the original Tacotron 2 has only been experimented
 182 with English. Therefore, instead of using the com-
 183 mon approach of Tacotron 2, we develop an end-to-
 184 end process using the Grad-TTS architecture³, a neu-
 185 ral network using the denoising diffusion probabilis-
 186 tic model. This approach is also consistent with re-
 187 lated studies in the group of languages closely related
 188 to Bahnaric, such as Vietnamese²⁴.
 189 Our text-to-speech system consists of three main
 190 components, as shown in Figure 5. First, the

191 *Text Analysis* module parses the text into a pseudo-
 192 phonetic representation, which is suitable for neural
 193 network processing.
 194 The second module is an acoustic model based on
 195 Grad-TTS, from the input of which is a set of pseudo-
 196 phonemes, it goes through the training process to
 197 generate the mel-spectrogram representation. Mel-
 198 spectrogram is a representation in the form of a spec-
 199 trum of sound waves, consisting of two dimensions,
 200 frequency and time. Mel-spectrograms can be ex-
 201 tracted directly from the sound wave and contain
 202 more detailed information about the frequency bands
 203 that prevail at each moment in the sound wave. Con-
 204 versely, it is also possible to extract sound waves from
 205 the mel-spectrogram through the inverse problem.

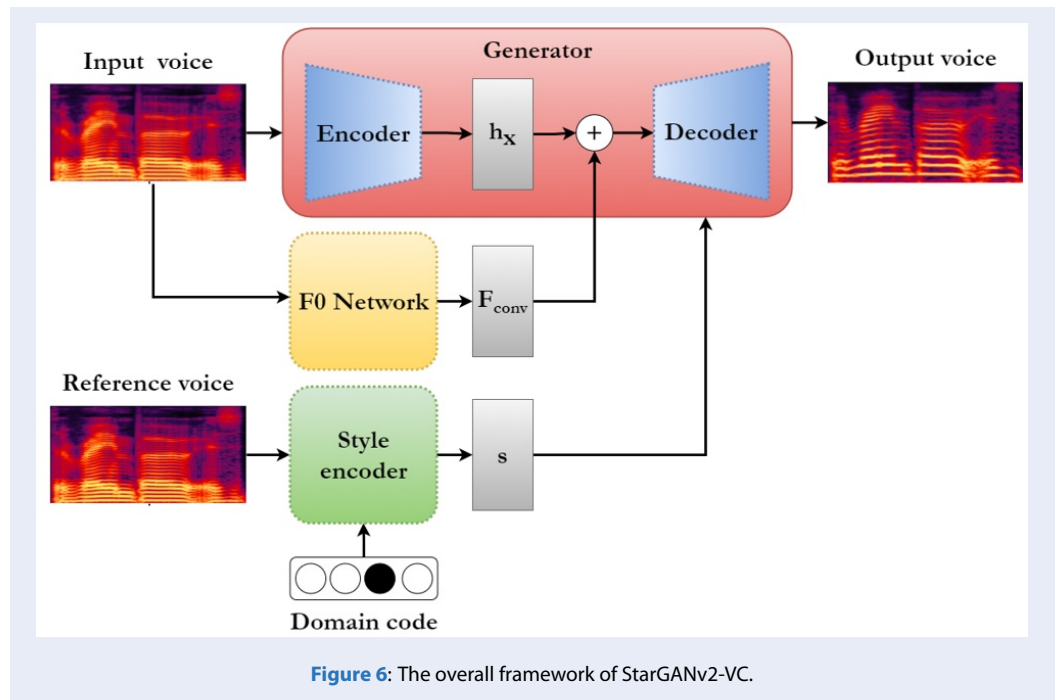


Figure 6: The overall framework of StarGANv2-VC.

206 The detailed architecture of this model can be con- 235
 207 sulted through related previous works, and in this 236
 208 publication, in order to be accessible to a wide audi- 237
 209 ence and not be too technical, we do not go through 238
 210 the details of this model. 239

211 The final step is performed by the vocoder. We use 240
 212 the HiFi-GAN network⁵ to convert the output from 241
 213 mel-spectrogram to waveform. More specifically, in- 242
 214 stead of using pre-trained HiFi-GAN for the English 243
 215 language, we retrained a pre-trained HiFi-GAN-BN 244
 216 system from Bahnaric to generate the final voice. 245

217 StarGANv2-VC model for Bahnaric voice 246 218 conversion 247

219 The Grad-TTS model can pronounce without lim- 248
 220 itation the vocabulary of Bahnar texts, but due to 249
 221 the low resource characteristics of this language, the 250
 222 sound quality still lacks the naturalness of humans. 251
 223 To overcome this problem, we propose to use the 252
 224 StarGANv2-VC model to convert the voice synthe- 253
 225 sized by Grad-TTS into a sample voice of the native 254
 226 Bahnar. The proposed methodology has been 255
 227 developed based on the foundational principles of 256
 228 StarGANv2-VC⁴, a pioneering framework that em- 257
 229 ploys a solitary discriminator and generator to pro- 258
 230 duce a diverse array of images across various domains. 259
 231 These domains are characterized by the utilization of 260
 232 domain-specific style vectors sourced either from the 261
 233 style encoder or the mapping network. In the do- 262
 234 main of voice conversion, each speaker is treated as a

235 discrete domain. To ensure the maintenance of con- 236
 237 sistent fundamental frequency (F0) conversion, the 238
 239 network architecture has been thoughtfully enhanced 240
 241 through the integration of a pre-trained joint detec- 242
 243 tion and classification (JDC) F0 extraction network²⁵. 244
 245 Figure 6, presented herein, offers an illustrative depic- 246
 247 tion of the StarGANv2-VC framework for elucidation. 248

249 In StarGANv2-VC, a sample $X \in X_{y_{src}}$ from the source 250
 251 domain $y_{src} \in Y$ undergoes transformation to a corre- 252
 253 sponding sample $\hat{X} \in X_{y_{trg}}$ in the target domain $y_{trg} \in 254$
 255 Y via a mapping function, denoted as $G : X_{y_{src}} \rightarrow X_{y_{trg}}$. 256
 257 Crucially, this transformation is achieved independ- 258
 259 ently of parallel data. 260

Throughout the training process, the selection of the 261
 262 target domain, $Y_{trg} \in Y$, is random, and its style code, 263
 264 $s \in S_{y_{trg}}$, is encoded through a style encoder. This en- 265
 266 coder utilizes a reference input $X_{ref} \in X$ from the tar- 267
 268 get domain to produce the style code, designated as 269
 270 $s = S(X_{ref}, y_{trg})$. Using a mel-spectrogram $X \in X_{y_{src}}$ 271
 272 from the source domain $y_{src} \in Y$ and the target do- 273
 274 main $y_{trg} \in Y$, our model is trained by minimizing the 275
 276 subsequent loss functions. 277

Adversarial loss. The generator is trained to produce 278
 279 a new mel-spectrogram, denoted as $G(X, s)$, from an 280
 281 input mel-spectrogram X and a style vector s by uti- 282
 283 lizing the adversarial loss. 284

$$L_{adv} = E_{X, y_{src}} [\log D(X, y_{src})] + E_{X, y_{trg}, s} [\log (1 - D(G(X, s), y_{trg}))] \quad (1)$$

261 where $D(\cdot, y)$ represents the output of the real/fake
 262 classifier of the domain $y \in Y$.
 263 **Adversarial source classifier loss.** Another adversar-
 264 ial loss function, involving the source classifier C , is
 265 employed (refer to Figure 7).

$$I_{advcls} = E_{X, y_{trg}, s} [CE(C(G(X, s), y_{trg}))] \quad (2)$$

266 where $CE(\cdot)$ denotes the cross-entropy loss function.
 267 **Style reconstruction loss.** To guarantee that the style
 268 code can be reconstructed from the generated sam-
 269 ples, the style reconstruction loss is used.

$$L_{sty} = E_{X, y_{trg}, s} [\|s - S(G(X, s), y_{trg})\|_1] \quad (3)$$

270 **Style diversification loss.** The different samples must
 271 be generated with different style codes. We enforce
 272 the generator to learn this constraint by maximizing
 273 the style diversification loss. In addition to maximiz-
 274 ing the mean absolute error (MAE) between gener-
 275 ated samples, the MAE of the F0 features between
 276 samples generated with different style codes is also
 277 maximized.

$$L_{ds} = E_{X, s_1, s_2, y_{trg}} [\|G(X, s_1) - G(X, s_2)\|_1] + \quad (4)$$

$$E_{X, s_1, s_2, y_{trg}} [\|F_{conv}(G(X, s_1)) - F_{conv}(G(X, s_2))\|_1]$$

278 where $s_1, s_2 \in S_{y_{trg}}$ are two randomly sampled style
 279 codes from domain $y_{trg} \in Y$ and $F_{conv}(\cdot)$ is the output
 280 of convolutional layers of F0 network F.

281 F0 consistency loss. An F0-consistent loss is added
 282 to produce F0-consistent results with the normalized
 283 F0 curve provided by F0 network F. For a given input
 284 mel-spectrogram X , the function $F(X)$ calculates the
 285 absolute fundamental frequency (F0) value in Hertz
 286 for each frame within X . Given that male and female
 287 speakers tend to exhibit distinct average F0 values, a
 288 normalization step is employed to standardize the ab-
 289 solute F0 values captured by $F(X)$. This normalization
 290 process is represented as $\hat{F}(X) = \frac{F(X)}{\|F(X)\|_1}$.
 291 Consequently, the F0 consistency loss is formulated as
 292 follows

$$L_{f0} = E_{X, s} [\|\hat{F}(X) - \hat{F}(G(X, s))\|_1] \quad (5)$$

293 **Speech consistency loss.** Ensuring the linguistic fi-
 294 delity of the converted speech is paramount, achieved
 295 through the implementation of a speech consistency
 296 loss mechanism. This mechanism relies on convo-
 297 lutional features extracted from a pre-trained joint
 298 Connectionist Temporal Classification (CTC) - at-
 299 tention model, particularly the VGG-Bidirectional
 300 Long Short-Term Memory (BLSTM) network, de-
 301 tailed in reference²⁶ and accessible within the Esp-
 302 net toolkit²⁷. Adhering to the approach of previous

research²⁸, we leverage the output from the interme-
 303 diate layer preceding the Long Short-Term Memory
 304 (LSTM) layers, denoted as $h_{asr}(\cdot)$, to encapsulate the
 305 linguistic feature. Consequently, the formal defini-
 306 tion of the speech consistency loss is as follows
 307

$$L_{asr} = E_{X, s} [\|h_{asr}(X) - h_{asr}(G(X, s))\|_1] \quad (6)$$

Norm consistency loss. In order to maintain the
 308 temporal integrity of generated samples, we employ
 309 a norm consistency loss. This loss mechanism is
 310 designed to ensure the preservation of speech and
 311 silence intervals in the generated output. To calcu-
 312 late the absolute column-sum norm for a mel-
 313 spectrogram X , which comprises N mel frequency
 314 bins and T frames at the t^{th} frame, we define it as
 315 $\|X_{:,t}\| = \sum_{n=1}^N \|X_{n,t}\|_1$, where $t \in \{1, \dots, T\}$ represents
 316 the frame index. The norm consistency loss can be
 317 expressed as follows
 318

$$L_{norm} = E_{X, s} \left[\frac{1}{T} \sum_1^T \|\|X_{:,t}\| - \|G(X, s)_{:,t}\|\| \right] \quad (7)$$

319 **Cycle consistency loss.** Finally, we introduce the
 320 cycle consistency loss, as outlined in reference¹⁷,
 321 with the purpose of preserving all remaining features
 322 present in the input data.

$$L_{cyc} = E_{X, y_{src}, y_{trg}, s} [\|X - G(G(X, s), \tilde{s})\|_1] \quad (8)$$

323 where $\tilde{s} = S(X, y_{src})$ is the estimated style code of the
 324 input in the source domain $y_{src} \in Y$.

325 **Full objective.** The entirety of our generator's objec-
 326 tive functions can be condensed as follows:

$$\begin{aligned} \min_{G, S, M} & L_{adv} + \lambda_{advcls} L_{advcls} + \lambda_{sty} L_{sty} \\ & - \lambda_{ds} L_{ds} + \lambda_{f0} L_{f0} + \lambda_{asr} L_{asr} \\ & + \lambda_{norm} L_{norm} + \lambda_{cyc} L_{cyc} \end{aligned} \quad (9)$$

327 where λ_{advcls} , λ_{sty} , λ_{ds} , λ_{f0} , λ_{asr} , λ_{norm} and λ_{cyc}
 328 are hyperparameters for each term.

329 The complete objective for our discriminator is as fol-
 330 lows

$$\min_{C, D} -L_{adv} + \lambda_{cls} L_{cls} \quad (10)$$

331 where λ_{cls} is the hyperparameter for source classifier
 332 loss L_{cls} , which is given by

$$L_{cls} = E_{X, y_{src}, s} [CE(C(G(X, s), y_{src}))] \quad (11)$$

333 **The pretrained HiFi-GAN-BN model from**
 334 **Bahnaric language for the vocoder of Grad-**
 335 **TTS model.**

336 Vocoders serve as instruments employed for trans-
 337 forming a speech spectrogram into audible sound
 338 waves. They play a pivotal role in the voice conversion
 339 process, facilitating the creation of sound correspond-
 340 ing to the given spectrogram. As outlined in Sec-
 341 tion 4.2, when it comes to the Grad-TTS system, em-
 342 ploying a pre-trained HiFi-GAN designed for the En-
 343 glish language poses several challenges due to the dis-
 344 tinct linguistic and acoustic characteristics inherent
 345 in the Bahnar language as opposed to English. Con-
 346 sequently, we took the approach of retraining a pre-
 347 existing HiFi-GAN system tailored to Bahnar voice,
 348 following the methodology illustrated in Figure 8.

349 Within this training pipeline, there are three key com-
 350 ponents: one generator and two discriminators. The
 351 generator, designed as a fully convolutional neural
 352 network, takes a mel-spectrogram as its input and em-
 353 ploys transposed convolutions to up-sample it until
 354 the resulting sequence matches the temporal resolu-
 355 tion of raw waveforms.

356 In terms of the discriminators, they consist of two dis-
 357 tinct modules. Firstly, the multi-period discrimina-
 358 tor (MPD) is composed of several sub-discriminators,
 359 each responsible for assessing specific segments of pe-
 360 riodic signals within the input audio. Furthermore,
 361 to capture consecutive patterns and long-term depen-
 362 dencies, we incorporate the multi-scale discriminator
 363 (MSD) concept, which is inspired by the approach in-
 364 troduced in MelGAN²⁹. This MSD evaluates audio
 365 samples at various levels to gain a comprehensive un-
 366 derstanding of the data.

367 The training process involves adversarial training for
 368 both the generator and discriminators. Additionally,
 369 two supplementary loss functions are employed to
 370 enhance training stability and overall model perfor-
 371 mance.

372 **GAN loss.** The training objectives of this model ad-
 373 here to the principles of LSGAN³⁰. Specifically, they
 374 replace the binary cross-entropy terms from the origi-
 375 nal GAN objectives³¹ with least squares loss functions
 376 to ensure non-vanishing gradient flows. In this setup,
 377 the discriminator’s training goal is to classify ground
 378 truth samples as 1 and generated samples from the
 379 generator as 0. Conversely, the generator aims to de-
 380 ceive the discriminator by adjusting the quality of its
 381 generated samples to be classified as a value very close
 382 to 1.

The GAN losses for both the generator G and the dis-
 383 criminator D are defined as 384

$$L_{adv}(D;G) = E_{X,s} \left[(D(X) - 1)^2 + (D(G(s)))^2 \right] \quad (12)$$

$$L_{adv}(G;D) = E_s \left[(D(G(s)) - 1)^2 \right] \quad (13)$$

where X denotes the ground truth audio and denotes the
 385 mel-spectrogram of the ground truth audio. 386

Mel-Spectrogram loss. To enhance the training per-
 387 formance of the generator and ensure the synthesized
 388 audio’s fidelity, we introduce a mel-spectrogram loss
 389 into the GAN objective. This addition is made with
 390 the expectation that the input condition should also
 391 play a role in improving the perceptual quality, tak-
 392 ing into consideration the characteristics of the hu-
 393 man auditory system. 394

The mel-spectrogram loss is calculated as the L1 dis-
 395 tance between the mel-spectrogram of a waveform
 396 generated by the generator and that of a ground truth
 397 waveform. It is defined as 398

$$L_{Mel}(G) = E_{X,s} [\|\phi(X) - \phi(G(s))\|_1] \quad (14)$$

where ϕ represents the transform function used to
 399 derive the mel-spectrogram from the corresponding
 400 waveform. 401

Feature matching loss. The model can also undergo
 402 optimization based on a metric that quantifies the
 403 distinction in features extracted by the discriminator
 404 when comparing a ground truth sample to a generated
 405 sample³². This metric, known as the feature matching
 406 loss, is defined as follows 407

$$L_{FM}(G;D) = E_{X,s} \left[\sum_{i=1}^T \frac{1}{N} \|D^i(X) - D^i(G(s))\|_1 \right] \quad (15)$$

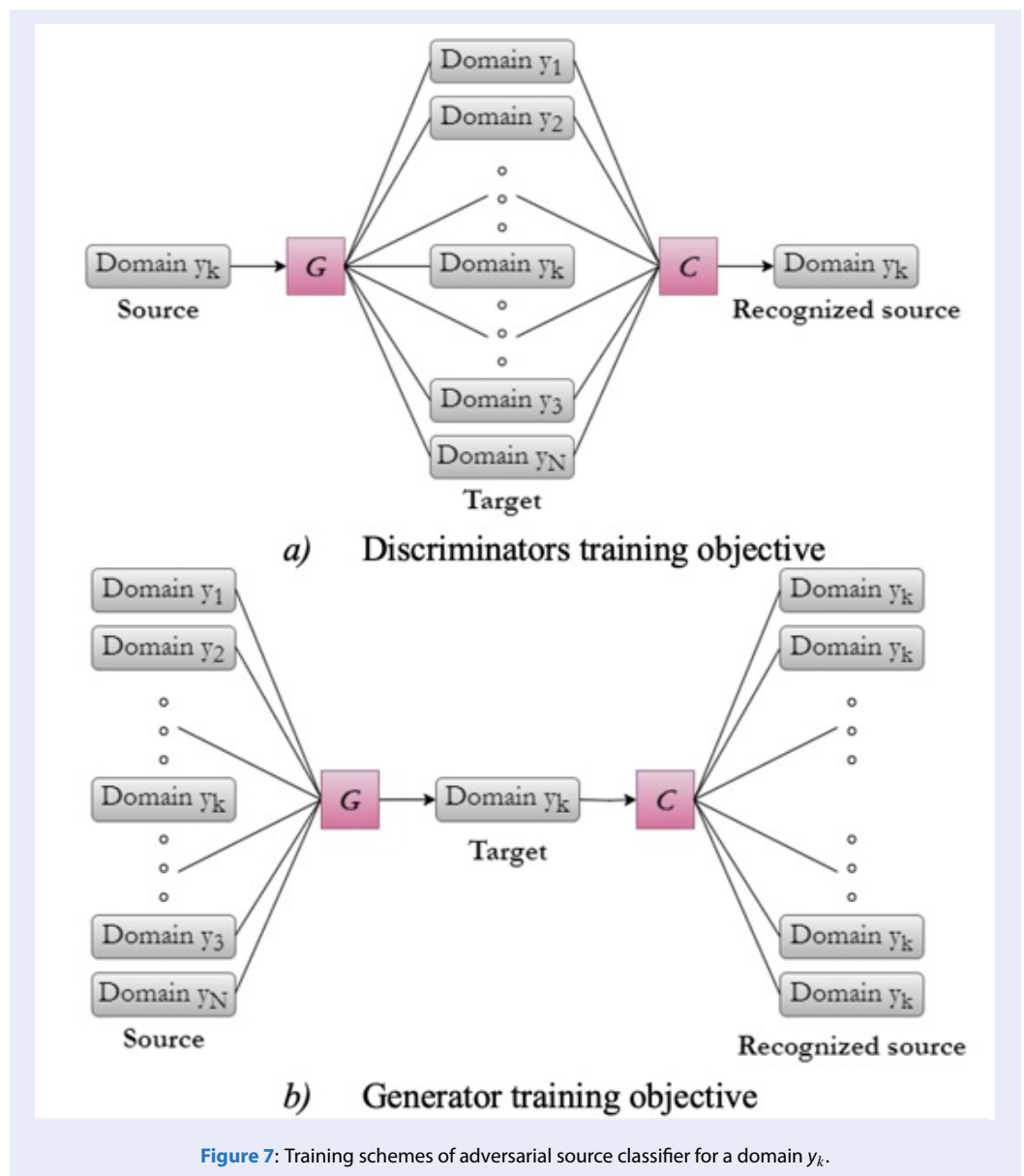
Full objective. The ultimate loss functions for both
 408 the generator and discriminator are defined as 409

$$\min_{G,D} L_G = L_{adv}(G,D) + \lambda_{fm} L_{FM}(G,D) + \lambda_{mel} L_{Mel}(G) \quad (16)$$

$$\min_{G,D} L_D = L_{adv}(D,G) \quad (17)$$

EXPERIMENT RESULTS 410

There are two main models trained from scratch in
 411 this system including the StarGANv2-VC model for
 412 voice conversion and the HiFi-GAN for the vocoder
 413 of the Grad-TTS model. Both two these models are
 414 developed based on the Pytorch framework. Consid-
 415 ering the StarGANv2-VC model, it is trained with 122
 416 epochs using the GPU of NVIDIA RTX 3080. The
 417 dataset that we use to train this model is the recorded
 418



419 voices gained manually by native Bahnaric from the
 420 provinces of Gia Lai, Kon Tum, and Binh Dinh in
 421 Vietnam, where exist considerable communities of
 422 Bahnaric people. They are used as training input for
 423 audio files that are generated from Grad-TTS. On the
 424 other hand, the HiFi-GAN model is trained up to 1
 425 million steps with two A100 GPUs. In other to train
 426 this model, we collected the from the YouTube chan-
 427 nel of VTV5, which consists of 300 hours of Bahnaric
 428 speech.

429 Regarding the evaluation methodology, we built the
 430 web application as shown in Figure 9. A user-friendly
 431 web-based interface was developed using Streamlit to

432 facilitate the evaluation process. This interface pre-
 433 sented users with 20 questions, each representing a
 434 unique evaluation instance. Each evaluation instance
 435 consisted of the following components:

436 **Original Speech Audio:** The interface played an origi-
 437 nal speech audio recording from a human speaker.
 438 This audio served as a reference point for users to
 439 compare the converted audio against.

440 **Converted Speech Audios:** Two converted speech
 441 audios were played for each evaluation instance.
 442 These audios were generated using our two best-
 443 performing StarGANv2-VC models. The intention
 444 here was to compare the quality of voice conversion
 445 between the models.

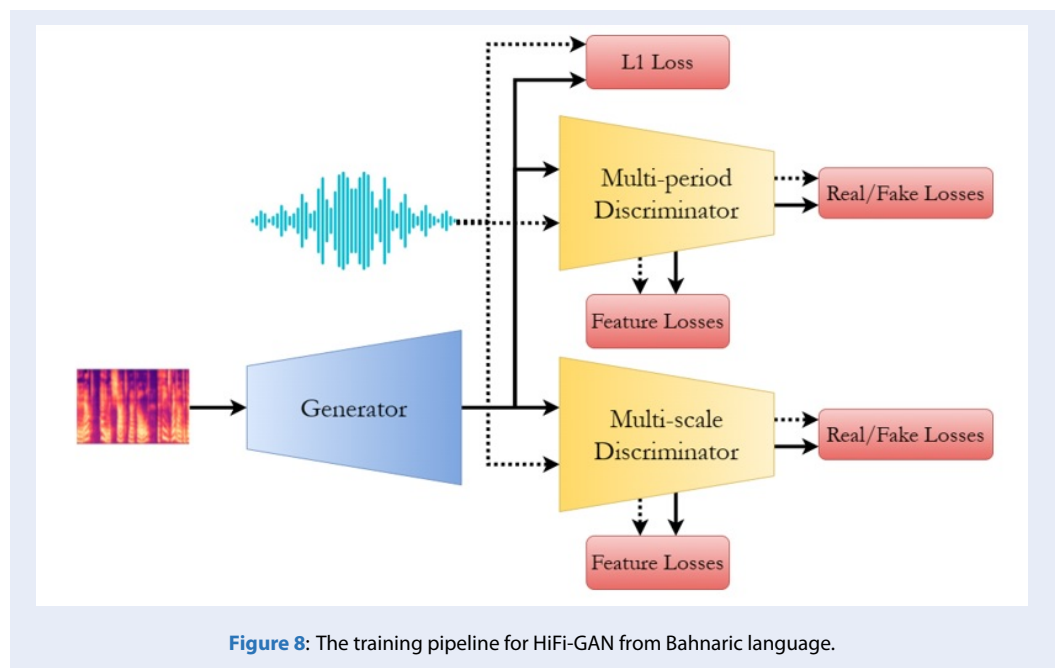


Figure 8: The training pipeline for HiFi-GAN from Bahharic language.

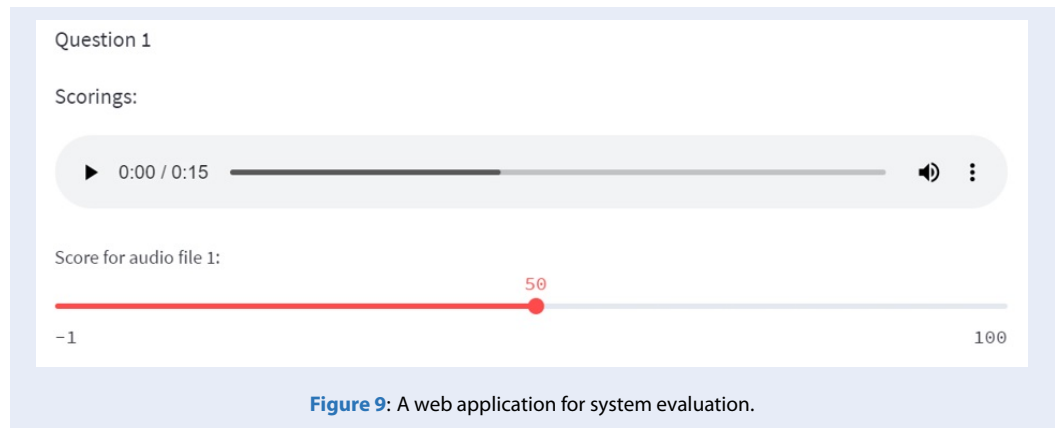


Figure 9: A web application for system evaluation.

446 With respect to the scoring mechanism, users were
 447 given a scoring scale ranging from -1 to 100 to rate
 448 the quality of the converted audio. This scoring scale
 449 was designed to capture a broad spectrum of quality
 450 perceptions. The interpretation of the scale was as fol-
 451 lows:

- 452 • **-1: Unrealistic Sound.** The converted audio
 453 needs to be more realistic and unconvincing to
 454 represent the target speaker.
- 455 • **0-49: Poor to Fair.** The converted audio is
 456 poor to fair quality, with significant discrepan-
 457 cies from the original speaker's voice.
- 458 • **50-69: Moderate.** The converted audio resem-
 459 bles the target speaker's voice, but improve-
 460 ments are needed.

- 461 • **70-89: Good.** The converted audio is of 461
 462 good quality and reasonably captures the target 462
 463 speaker's characteristics. 463
- 464 • **90-99: Very Good.** The converted audio is of 464
 465 outstanding quality, closely resembling the tar- 465
 466 get speaker's voice with minor discrepancies. 466
- 467 • **100: Perfect.** The converted audio is indistin- 467
 468 guishable from the audio of the actual target 468
 469 speaker; no improvements are necessary. 469

The scale ranging from -1 to 100 (comprising 6 lev- 470
 471 els) has been designed with specific intentions. At the 471
 472 lower end, -1 is assigned to instances where the AI- 472
 473 generated sound is exceptionally poor, to the extent 473
 474 that it is practically unbearable. Conversely, at the 474
 475 upper end, 100 signifies that within the provided au- 475
 476 dio, at least one sentence closely resembles an original 476

477 human-generated recording. Essentially, the scale is
 478 employed to convey to the evaluator that there can be
 479 a wide range of sound quality, spanning from severely
 480 subpar to human-level excellence. The evaluation re-
 481 sult is collected from 46 voluntary participants, whose
 482 statistics are shown in Table 1.

483 As shown in Table 1, there is no evaluation result
 484 of bad quality in the samples of the original voice
 485 that recorded by native speakers. Concerning voice
 486 conversion models, the VC-original model is trained
 487 from original voice data and the VC-Grad-TTS is
 488 trained with a suitable amount of data in the source
 489 domain that is taken from the output of Grad-TTS.

490 It can be seen that the VC-original model generates
 491 sounds with acceptable quality. However, there is an
 492 existence of bad quality samples and it accounts for
 493 4.24% of the evaluation set. The number of samples
 494 having very good quality is also quite low at 11.96%.
 495 Overall, the voice converted by this model is evaluated
 496 as having good quality with a mean score of 74.07.

497 On the other hand, the VC-Grad-TTS model gives
 498 better performance. The number of samples that have
 499 poor to fair quality is reduced significantly (account-
 500 ing for 0.87%). In addition, most generated sample
 501 from this model is evaluated from good to perfect.
 502 The mean evaluation score is also high with 80.33,
 503 which belongs to the scale of good quality sound.

504 **DISCUSSION**

505 This research addresses the challenge of generating
 506 natural-sounding speech in the Bahnaric language,
 507 which is often marginalized and lacks adequate re-
 508 sources. Our system shows promising results in syn-
 509 thesizing Bahnaric speech. Table 1 illustrates that
 510 models trained with Grad-TTS output as the domain
 511 source outperform those trained directly with na-
 512 tive speaker data, with synthesized voice quality also
 513 rated as good. Moreover, the HiFi-GAN-BN model,
 514 pre-trained with Bahnaric voice data, enhances the
 515 authenticity of synthesized speech to resemble Bah-
 516 naric accents when converting mel-spectrogram out-
 517 put. On the other hand, further optimization and
 518 evaluation across diverse linguistic and cultural con-
 519 texts are necessary. Collaboration with linguists and
 520 community stakeholders is vital to ensure the cul-
 521 tural relevance and acceptance of synthesized Bah-
 522 naric voices. Ultimately, our work contributes to the
 523 preservation and promotion of cultural diversity and
 524 linguistic heritage, not only within the Bahnaric com-
 525 munity in Vietnam but also in similar contexts world-
 526 wide.

CONCLUSION

The Vietnamese government is endeavoring to en-
 enhance their integration through advancements in
 socio-cultural and scientific literacy. In order to
 contribute to conveying information with native-
 like Bahnaric speech, we have proposed an effec-
 tive approach called BN-TTS-VC system. Most
 of the text-to-speech systems require a substantial
 amount of training data. It is particularly ardu-
 ous to gather a high-quality training dataset of ex-
 tremely low-resource languages like Bahnaric. There-
 fore, our system combined Grad-TTS model³ and the
 StarGANv2-VC model⁴ to solve this problem. In ad-
 dition, we also introduce the HiFi-GAN-BN model, a
 variant of HiFi-GAN⁵ pre-trained by Bahnaric voice,
 to resemble the Bahnaric accents better when trans-
 forming the mel-spectrogram output of StarGANv2-
 VC into human-listenable waveform. The evaluation
 results have shown that the system is able to generate
 good-quality audio and the voice conversion model
 that is trained with the source domain data taken from
 the output of Grad-TTS gives better performance. Fu-
 ture work includes improving the quality of sound
 that is not clear or missing the vocabulary of the text.

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 sys-
 tem (all dialects)” - KC-4.0-29/19-25

LIST OF ABBREVIATIONS

TTS: Text-to-speech
 VC: Voice conversion
 Grad-TTS: A Diffusion Probabilistic Model for Text-
 to-Speech
 StarGANv2-VC: A Diverse, Unsupervised, Non-
 parallel Framework for Natural-Sounding Voice Con-
 version.
 HiFi-GAN: A GAN-based model capable of generat-
 ing high fidelity speech efficiently.
 BN-TTS-VC: The combined system of text-to-speech
 and voice conversion for Bahnaric language.
 HiFi-GAN-BN: A GAN-based model from Bahnaric
 language for the vocoder of Grad-TTS model.

CONFLICTS OF INTEREST

All authors declare that they have no conflicts of in-
 terest.

Table 1: The evaluation result of StarGANv2-VC models.

Type of sample	Quality (%)						Mean score
	-1 ↓	0-49 ↓	50-69 ↑	70-89 ↑	90-99 ↑	100 ↑	
Original	0.0	0.0	2.06	56.31	39.02	2.61	87.12
VC-original	0.0	4.24	30.22	52.39	11.96	1.19	74.07
VC-Grad-TTS	0.0	0.87	18.26	55.54	23.59	1.74	80.33

CREDIT AUTHORSHIP CONTRIBUTION STATEMENT

576 **Dang Tran Dat:** Methodology, Model development,
 577 Evaluation, Writing – Original Draft.
 578 **Tang Quoc Thai:** Methodology, Model development,
 579 Evaluation, Writing.
 580 **Nguyen Quang Duc:** Methodology, System De-
 581 ployment, Resources, Data Collection, Data Curation,
 582 Writing.
 583 **Vo Duy Hung:** Methodology.
 584 **Quan Thanh Tho:** Supervision, Project Administra-
 585 tion, Methodology, Writing - Review & Editing.

REFERENCES

588 1. Tan X, Chen J, Liu H, Cong J, Zhang C, Liu Y, Wang X,
 589 Leng Y, Yi Y, He L, et al. *Naturalspeech: End-to-end text*
 590 *to speech synthesis with human-level quality.* arXiv preprint
 591 arXiv:2205.04421; 2022;
 592 2. Sisman B, Yamagishi J, King S, Li H. An overview of voice
 593 conversion and its challenges: From statistical modeling to
 594 deep learning. *IEEE/ACM Transactions on Audio, Speech, and*
 595 *Language Processing.* 2020;29:132-157; Available from: <https://doi.org/10.1109/TASLP.2020.3038524>.
 596 3. Popov V, Vovk I, Gogoryan V, Sadekova T, Kudinov M. *Grad-tts:*
 597 *A diffusion probabilistic model for text-to-speech.* In: *International*
 598 *Conference on Machine Learning;* 2021;.
 599 4. Choi Y, Uh Y, Yoo J, Ha J-W. *Stargan v2: Diverse image*
 600 *synthesis for multiple domains.* In: *Proceedings of*
 601 *the IEEE/CVF conference on computer vision and pattern*
 602 *recognition;* 2020; Available from: [https://doi.org/10.1109/](https://doi.org/10.1109/CVPR42600.2020.00821)
 603 [CVPR42600.2020.00821](https://doi.org/10.1109/CVPR42600.2020.00821).
 604 5. Kong J, Kim J, Bae J. *Hifi-gan: Generative adversarial networks*
 605 *for efficient and high fidelity speech synthesis.* In: *Advances*
 606 *in Neural Information Processing Systems.* 2020;33:17022-
 607 17033;.
 608 6. Scully C. *Articulatory synthesis.* In: *Speech production and*
 609 *speech modelling.* Springer; 1990. p. 151-186; Available from:
 610 https://doi.org/10.1007/978-94-009-2037-8_7.
 611 7. Lukose S, Upadhya SS. *Text to speech synthesizer-formant*
 612 *synthesis.* In: *2017 International Conference on Nascent Tech-*
 613 *nologies in Engineering (ICNTE);* 2017; PMID: 29031741. Avail-
 614 able from: <https://doi.org/10.1109/ICNTE.2017.7947945>.
 615 8. Mao X, Li Q, Xie H, Lau RYK, Wang Z, Smolley SP. *Least squares*
 616 *generative adversarial networks.* In: *Proceedings of the IEEE*
 617 *international conference on computer vision;* 2017; Available
 618 from: <https://doi.org/10.1109/ICCV.2017.304>.
 619 9. Kumar K, Kumar R, De Boissiere T, Gestin L, Teoh WZ, Sotelo
 620 J, De Brebisson A, Bengio Y, Courville AC. *Melgan: Generative*
 621 *adversarial networks for conditional waveform synthesis.* In:
 622 *Advances in neural information processing systems.* 2019;32;
 623 10. Park J, Zhao K, Peng K, Ping W. *Multi-speaker end-to-end*
 624 *speech synthesis.* arXiv preprint arXiv:1907.04462; 2019;.
 625 11. Polyak A, Wolf L, Adi Y, Taigman Y. *Unsupervised cross-domain*
 626 *singing voice conversion.* arXiv preprint arXiv:2008.02830;

2020; Available from: <https://doi.org/10.21437/Interspeech.2020-1862>.
 627 12. Watanabe S, Hori T, Karita S, Hayashi T, Nishitoba J, Unno
 628 Y, Soplin NEY, Heymann J, Wiesner M, Chen N, et al. *Es-*
 629 *pnet: End-to-end speech processing toolkit.* arXiv preprint
 630 arXiv:1804.00015; 2018; PMID: 29730221. Available from:
 631 <https://doi.org/10.21437/Interspeech.2018-1456>.
 632 13. Kim S, Hori T, Watanabe S. *Joint CTC-attention based end-to-*
 633 *end speech recognition using multi-task learning.* In: *2017*
 634 *IEEE international conference on acoustics, speech and signal*
 635 *processing (ICASSP);* 2017; Available from: [https://doi.org/10.](https://doi.org/10.1109/ICASSP.2017.7953075)
 636 [1109/ICASSP.2017.7953075](https://doi.org/10.1109/ICASSP.2017.7953075).
 637 14. Kum S, Nam J. *Joint detection and classification of singing*
 638 *voice melody using convolutional recurrent neural networks.*
 639 *Applied Sciences.* 2019;9:1324; Available from: [https://doi.org/](https://doi.org/10.3390/app9071324)
 640 [10.3390/app9071324](https://doi.org/10.3390/app9071324).
 641 15. Tran T, Nguyen T, Bui H, Nguyen K, Vo NG, Pham TV, Quan
 642 T. *Naturalness Improvement of Vietnamese Text-to-Speech*
 643 *System Using Diffusion Probabilistic Modelling and Unsuper-*
 644 *vised Data Enrichment.* In: *International Conference on Intelli-*
 645 *gence of Things;* 2022; Available from: [https://doi.org/10.1007/](https://doi.org/10.1007/978-3-031-15063-0_36)
 646 [978-3-031-15063-0_36](https://doi.org/10.1007/978-3-031-15063-0_36).
 647 16. Huang W-C, Hayashi T, Wu Y-C, Kameoka H, Toda T. *Voice*
 648 *transformer network: Sequence-to-sequence voice conver-*
 649 *sion using transformer with text-to-speech pretraining.* arXiv
 650 preprint arXiv:1912.06813; 2019; Available from: [https://doi.](https://doi.org/10.21437/Interspeech.2020-1066)
 651 [org/10.21437/Interspeech.2020-1066](https://doi.org/10.21437/Interspeech.2020-1066).
 652 17. Tanaka K, Kameoka H, Kaneko T, Hojo N. *AttS2S-VC: Sequence-*
 653 *to-sequence voice conversion with attention and context*
 654 *preservation mechanisms.* In: *ICASSP 2019-2019 IEEE Interna-*
 655 *tional Conference on Acoustics, Speech and Signal Process-*
 656 *ing (ICASSP);* 2019; Available from: [https://doi.org/10.1109/](https://doi.org/10.1109/ICASSP.2019.8683282)
 657 [ICASSP.2019.8683282](https://doi.org/10.1109/ICASSP.2019.8683282).
 658 18. Zhu J-Y, Park T, Isola P, Efros AA. *Unpaired image-to-image*
 659 *translation using cycle-consistent adversarial networks.* In:
 660 *Proceedings of the IEEE international conference on com-*
 661 *puter vision;* 2017; Available from: [https://doi.org/10.1109/](https://doi.org/10.1109/ICCV.2017.244)
 662 [ICCV.2017.244](https://doi.org/10.1109/ICCV.2017.244).
 663 19. Park S-w, Kim D-y, Joe M-c. *Cotatron: Transcription-guided*
 664 *speech encoder for any-to-many voice conversion without*
 665 *parallel data.* arXiv preprint arXiv:2005.03295; 2020; Available
 666 from: <https://doi.org/10.21437/Interspeech.2020-1542>.
 667 20. Kaneko T, Kameoka H, Tanaka K, Hojo N. *Cyclegan-vc3: Exam-*
 668 *ining and improving cyclegan-vc3 for mel-spectrogram*
 669 *conversion.* arXiv preprint arXiv:2010.11672; 2020; Available from:
 670 <https://doi.org/10.21437/Interspeech.2020-2280>.
 671 21. Huang W-C, Luo H, Hwang H-T, Lo C-C, Peng Y-H, Tsao
 672 Y, Wang H-M. *Unsupervised representation disentanglement*
 673 *using cross domain features and adversarial learning in vari-*
 674 *ational autoencoder based voice conversion.* *IEEE*
 675 *Transactions on Emerging Topics in Computational Intelli-*
 676 *gence.* 2020;4:468-479; Available from: [https://doi.org/10.](https://doi.org/10.1109/TETCI.2020.2977678)
 677 [1109/TETCI.2020.2977678](https://doi.org/10.1109/TETCI.2020.2977678).
 678 22. Ding S, Gutierrez-Osuna R. *Group Latent Embedding for Vec-*
 679 *tor Quantized Variational Autoencoder in Non-Parallel Voice*
 680 *Conversion.* In: *Interspeech;* 2019; PMID: 31791587. Available
 681 from: <https://doi.org/10.21437/Interspeech.2019-1198>.
 682 23. Qian K, Zhang Y, Chang S, Yang X, Hasegawa-Johnson M. Au-

- 687 tovc: Zero-shot voice style transfer with only autoencoder
688 loss. In: International Conference on Machine Learning; 2019;.
689 24. Prenger R, Valle R, Catanzaro B. Waveglow: A flow-based generative
690 network for speech synthesis. In: ICASSP 2019-2019
691 IEEE International Conference on Acoustics, Speech and Signal
692 Processing (ICASSP); 2019; Available from: [https://doi.org/
693 10.1109/ICASSP.2019.8683143](https://doi.org/10.1109/ICASSP.2019.8683143).
- 694 25. Lam QT, Do DH, Vo TH, Nguyen DD. Alternative vietnamese
695 speech synthesis system with phoneme structure. In: 2019
696 19th International Symposium on Communications and Information
697 Technologies (ISCIT); 2019; Available from: [https://doi.
698 org/10.1109/ISCIT.2019.8905142](https://doi.org/10.1109/ISCIT.2019.8905142).
- 699 26. Shen J, Pang R, Weiss RJ, Schuster M, Jaitly N, Yang Z, Chen
700 Z, Zhang Y, Wang Y, Skerrv-Ryan R, et al. Natural tts synthesis
701 by conditioning wavenet on mel spectrogram predictions.
702 In: 2018 IEEE international conference on acoustics, speech
703 and signal processing (ICASSP); 2018; Available from: [https:
704 //doi.org/10.1109/ICASSP.2018.8461368](https://doi.org/10.1109/ICASSP.2018.8461368).
- 705 27. Donahue J, Dieleman S, Bińkowski M, Elsen E, Simonyan
706 K. End-to-end adversarial text-to-speech. arXiv preprint
707 arXiv:2006.03575; 2020;.
- 708 28. Ren Y, Hu C, Tan X, Qin T, Zhao S, Zhao Z, Liu T-Y. Fastspeech
709 2: Fast and high-quality end-to-end text to speech. arXiv
710 preprint arXiv:2006.04558; 2020;.
- 711 29. Zen H, Tokuda K, Black AW. Statistical parametric speech syn-
712 thesis. Speech Communication. 2009;51:1039-1064; Available
713 from: <https://doi.org/10.1016/j.specom.2009.04.004>.
- 714 30. Schwarz D. Corpus-based concatenative synthesis. IEEE signal
715 processing magazine. 2007;24:92-104; Available from: [https://
716 doi.org/10.1109/MSP.2007.323274](https://doi.org/10.1109/MSP.2007.323274).
- 717 31. Goodfellow I, Pouget-Abadie J, Mirza M, Xu B, Warde-Farley
718 D, Ozair S, Courville A, Bengio Y. Generative adversarial nets.
719 Advances in neural information processing systems. 2014;27;.
- 720 32. Larsen ABL, Sønderby SK, Larochelle H, Winther O. Autoen-
721 coding beyond pixels using a learned similarity metric. In: In-
722 ternational conference on machine learning; 2016;.

Phương pháp thay đổi giọng tăng cường tính tự nhiên cho quá trình sinh giọng nói ở ngôn ngữ ít tài nguyên: Thí nghiệm với ngôn ngữ Ba Na

Đặng Trần Đạt^{1,2}, Tăng Quốc Thái^{1,2}, Nguyễn Quang Đức^{1,2}, Võ Duy Hùng^{1,2}, Quán Thành Thơ^{1,2,*}



Use your smartphone to scan this QR code and download this article

¹Khoa Khoa học và Kỹ thuật Máy tính, Trường Đại học Bách khoa – ĐHQG-HCM, Việt Nam

²Đại học Quốc gia Thành phố Hồ Chí Minh, Việt Nam

Liên hệ

Quán Thành Thơ, Khoa Khoa học và Kỹ thuật Máy tính, Trường Đại học Bách khoa – ĐHQG-HCM, Việt Nam

Đại học Quốc gia Thành phố Hồ Chí Minh, Việt Nam

Email: qtttho@hcmut.edu.vn

Lịch sử

- Ngày nhận: 08-9-2023
- Ngày chấp nhận: 27-3-2024
- Ngày đăng:

DOI:



Bản quyền

© ĐHQG Tp.HCM. Đây là bài báo công bố mở được phát hành theo các điều khoản của the Creative Commons Attribution 4.0 International license.



TÓM TẮT

Ba Na là một nhóm dân tộc thiểu số ở Việt Nam, được chính phủ ưu tiên bảo tồn di sản văn hóa, truyền thống và ngôn ngữ. Trong kỷ nguyên của công nghệ AI hiện nay, việc tổng hợp giọng nói tiếng Ba Na để hỗ trợ những nỗ lực bảo tồn này chứa đựng tiềm năng đáng kể. Mặc dù công nghệ chuyển đổi giọng nói đã có những bước tiến trong việc nâng cao chất lượng và tính tự nhiên của giọng nói được tổng hợp nhưng nó chỉ được chú trọng phát triển chủ yếu đối với các ngôn ngữ được sử dụng rộng rãi. Do đó, các ngôn ngữ có nguồn tài nguyên hạn chế như ngôn ngữ thuộc họ tiếng Ba Na gặp nhiều khó khăn trong việc tổng hợp giọng nói. Nghiên cứu này giải quyết thách thức lớn trong việc tổng hợp giọng nói có tính tự nhiên ở các ngôn ngữ có nguồn tài nguyên thấp bằng cách khám phá các ứng dụng của kỹ thuật chuyển đổi giọng nói cho tiếng Ba Na. Chúng tôi giới thiệu hệ thống BN-TTS-VC, một phương pháp tiên phong tích hợp hệ thống chuyển văn bản thành giọng nói dựa trên Grad-TTS, với các kỹ thuật chuyển đổi giọng nói dựa trên StarGANv2-VC, và cả hai đều được thiết kế riêng cho các sắc thái của tiếng Ba Na. Grad-TTS cho phép hệ thống phát âm các từ trong ngôn ngữ Ba Na mà không bị giới hạn từ vựng, trong khi StarGANv2-VC nâng cao tính tự nhiên của giọng nói được tổng hợp, đặc biệt là trong bối cảnh các ngôn ngữ có nguồn tài nguyên thấp như tiếng Ba Na. Ngoài ra, chúng tôi còn giới thiệu mô hình HiFi-GAN được tinh chỉnh bằng tiếng Ba Na để nâng cao chất lượng giọng nói so với giọng bản địa, đảm bảo thể hiện giọng nói tiếng Ba Na chân thực hơn. Để đánh giá hiệu quả của phương pháp tiếp cận, chúng tôi đã tiến hành thử nghiệm dựa trên đánh giá của con người từ các tình nguyện viên. Các kết quả sơ bộ đầy hứa hẹn, cho thấy phương pháp của chúng tôi chứa nhiều tiềm năng trong việc tổng hợp giọng nói mang tính tự nhiên tiếng Ba Na. Qua nghiên cứu này, mục tiêu của chúng tôi là đóng góp vào các nỗ lực để bảo tồn và thúc đẩy di sản ngôn ngữ và văn hóa của nhóm dân tộc thiểu số Bahnar. Bằng cách tận dụng sức mạnh của công nghệ AI, chúng tôi mong muốn thu hẹp khoảng cách trong tổng hợp giọng nói cho các ngôn ngữ nguồn tài nguyên thấp và tạo điều kiện thuận lợi cho việc bảo tồn di sản văn hóa quý báu của họ.

Từ khóa: Tổng hợp giọng nói tiếng Ba Na, chuyển văn bản thành giọng nói, chuyển đổi giọng nói tự nhiên

Trích dẫn bài báo này: Đạt D T, Thái T Q, Đức N Q, Hùng V D, Thơ Q T. Phương pháp thay đổi giọng tăng cường tính tự nhiên cho quá trình sinh giọng nói ở ngôn ngữ ít tài nguyên: Thí nghiệm với ngôn ngữ Ba Na. *Sci. Tech. Dev. J. - Eng. Tech.* 2024; ():1-1.