

# Accent-VITS: accent transfer for end-to-end TTS

Linhan Ma<sup>1</sup>, Yongmao Zhang<sup>1</sup>, Xinfu Zhu<sup>1</sup>, Yi Lei<sup>1</sup>, Ziqian Ning<sup>1</sup>, Pengcheng Zhu<sup>2</sup>, and Lei Xie<sup>1\*</sup>

<sup>1</sup> Audio, Speech and Language Processing Group (ASLP@NPU), School of Computer Science, Northwestern Polytechnical University, Xi'an, China

<sup>2</sup> Fuxi AI Lab, NetEase Inc., Hangzhou, China

**Abstract.** Accent transfer aims to transfer an accent from a source speaker to synthetic speech in the target speaker's voice. The main challenge is how to effectively disentangle speaker timbre and accent which are entangled in speech. This paper presents a VITS-based [7] end-to-end accent transfer model named *Accent-VITS*. Based on the main structure of VITS, Accent-VITS makes substantial improvements to enable effective and stable accent transfer. We leverage a hierarchical CVAE structure to model accent pronunciation information and acoustic features, respectively, using bottleneck features and mel spectrums as constraints. Moreover, the text-to-wave mapping in VITS is decomposed into text-to-accent and accent-to-wave mappings in Accent-VITS. In this way, the disentanglement of accent and speaker timbre becomes more stable and effective. Experiments on multi-accent and Mandarin datasets show that Accent-VITS achieves higher speaker similarity, accent similarity and speech naturalness as compared with a strong baseline<sup>3</sup>.

**Keywords:** Text to speech · Accent transfer · Variational autoencoder · Hierarchical.

## 1 Introduction

In recent years, there have been significant advancements in neural text-to-speech (TTS), which can generate human-like natural speech from input text. Accented speech is highly desired for a better user experience in many TTS applications. Cross-speaker accent transfer is a promising technology for accented speech synthesis, which aims to transfer an accent from a source speaker to the synthetic speech in the target speaker's voice. Accent transfer can promote cross-region communication and make a TTS system better adapt to diverse language environments and user needs.

An accent is usually reflected in the phoneme pronunciation pattern and prosody variations, both of which are key attributes of the accent rendering [12, 14, 15]. The segmental and suprasegmental structures may be in distinctive pronunciation patterns for different accents and influence the listening perception

\* Corresponding author.

<sup>3</sup> Demos: <https://anonymous-accentvits.github.io/AccentVITS/>

of speaking accents [8, 25]. The prosody variations of accent are characterized by different pitch, energy, duration, and other prosodic appearance. To build an accent TTS system using accent transfer, the research problem can be treated as how to effectively *disentangle* speaker timbre and accent factors in speech. However, it is difficult to force the system to sufficiently disentangle the accent from the speaker timbre and content in speech since both pronunciation and prosody attributes are featured by local variations at the fine-grained level. And usually, each speaker has only one accent in the training phase which adds to the difficulty of disentangling.

Previous approaches attempting to disentangle accent attributes and speaker timbre are mainly based on Domain Adversarial Training (DAT) [5]. However, when the feature extraction function has a high capacity, DAT poses a weak constraint to the feature extraction function. Therefore, a single classifier with a gradient reversal layer in accent transfer TTS cannot disentangle the accent from the speaker’s timbre, as the accent is varied in prosody and pronunciation. Additionally, gradient descent in domain adversarial training can violate the optimizer’s asymptotic convergence guarantees, often hindering the transfer performance [1]. Applying DAT in accent transfer tasks, especially when each speaker has only one accent in the training phase, may result in inefficient and unstable feature disentanglement. Furthermore, there is a trade-off between speaker similarity and accent similarity, which means entirely removing speaker timbre hurts performance on preserving accent pronunciation [18].

Bottleneck (BN) features are recently used as an intermediate representation to supervise accent attribute modeling in TTS [24]. The BN feature, extracted from a well-trained neural ASR model, is considered to be noise-robust and speaker-independent [11, 19], which benefits speaker timbre and accent disentanglement. However, in the methods with BN as an intermediate representation [2, 13], models are often trained independently in multiple stages. This can lead to the issue of error accumulation and model mismatch between each stage, resulting in the degradation of synthesized speech quality and accent attributes.

In this paper, we propose an end-to-end accent transfer model, *Accent-VITS*, with a hierarchical conditional variational autoencoder (CVAE) [10] utilizing bottleneck features as a constraint to eliminate speaker timbre from the original signal. Specifically, we leverage the end-to-end speech synthesis framework, VITS [7], as the backbone of our model, since it achieves good audio quality and alleviates the error accumulation caused by the conventional two-stage TTS system consisting of an acoustic model and a vocoder. Based on the VITS structure, an additional CVAE is added to extract an accent-dependent latent distribution from the BN feature. The latent representation contains the accent and linguistic content and is modeled by the accented phoneme sequence input. The BN constraint factorizes the cross-speaker accent TTS into two joint-training processes, which are text-to-accent and accent-to-wave. The *text-to-accent* process takes the accented phoneme sequence as input to generate an accent-dependent distribution. The *accent-to-wave* process produces the speech distribution in the target accent and target speaker from the output accent distribution and is

conditioned on speaker identity. This design enables more effective learning of accent attributes, leading to sufficient disentanglement and superior performance of accent transfer for the synthesized speech. Experimental results on Mandarin multi-accent datasets demonstrate the superiority of our proposed model.

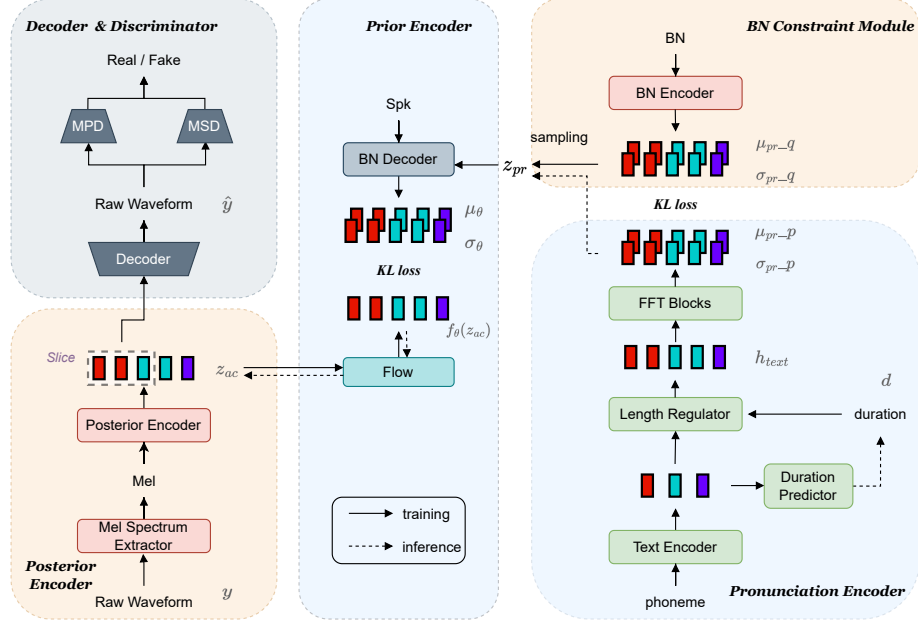


Fig. 1. Overview of Accent-VITS structure.

## 2 Method

This paper proposes a VITS-based end-to-end accent transfer model with a hierarchical conditional variational autoencoder (CVAE). As shown in Fig. 1, it mainly includes five parts: a posterior encoder, a decoder, a prior encoder, a pronunciation encoder, and a BN constraint module. The posterior encoder extracts the latent representation  $z_{ac}$  of acoustic feature from the waveform  $y$ , and the decoder reconstructs the waveform  $\hat{y}$  according to  $z_{ac}$ :

$$z_{ac} = \text{PostEnc}(y) \sim q(z_{ac}|y) \quad (1)$$

$$\hat{y} = \text{Dec}(z_{ac}) \sim p(y|z_{ac}) \quad (2)$$

The prior encoder produces a prior distribution of  $z_{ac}$ . Since the prior encoder of CVAE in VITS cannot effectively disentangle accent from the speaker timbre and text content, we use a hierarchical CVAE structure to model accent

information and acoustic features sequentially. We take the bottleneck feature (BN) that is extracted from the source wav by an ASR system as a constraint to improve the pronunciation information of accent in the latent space. The BN Encoder extracts the latent representation  $z_{pr}$  which contains accent pronunciation information from BN. The pronunciation encoder gets a prior distribution  $p(z_{pr}|c)$  of the latent variables  $z_{pr}$  given accented phoneme sequence condition  $c$ . The BN decoder in the prior encoder module also acts as the decoder of the first CVAE structure, getting the prior distribution  $p(z_{ac}|z_{pr}, spk)$  of the latent variables  $z_{ac}$  from the sampled latent representation  $z_{pr}$  given speaker identity condition  $spk$ . This hierarchical CVAE adopts a reconstruction objective  $L_{recon}$  and two prior regularization terms as

$$L_{cvae} = \alpha L_{recon} + D_{KL}(q(z_{pr}|BN)||p(z_{pr}|c)) + D_{KL}(q(z_{ac}|y)||p(z_{ac}|z_{pr}, spk)) \quad (3)$$

where  $D_{KL}$  is the Kullback-Leibler divergence. For the reconstruction loss, we use L1 distance of the mel-spectrum between ground truth and generated waveform. In the following, we will introduce the details of these modules.

## 2.1 Pronunciation Encoder

We assign a different phoneme set to each accent (standard Mandarin or accent Mandarin in this paper) and use a rule-based converter (G2P tool) to get the phoneme sequence from the text. Given the phoneme sequence of accent or Mandarin condition  $c$ , the pronunciation encoder module predicts the prior distribution  $p(z_{pr}|c)$  used for the prior regularization term of the first CVAE structure. In this module, the text encoder which consists of multiple FFT [20] blocks takes phoneme sequences as input and produces phoneme-level representation. Different from VITS, we use the length regulator (LR) in FastSpeech [16] to extend the phoneme-level representation to frame-level representation  $h_{text}$  [23]. The other multiple FFT blocks are used to extract a sequence of hidden vectors from the frame-level representation  $h_{text}$  and then generate the mean  $\mu_{pr\_p}$  and variance  $\sigma_{pr\_p}$  of the prior normal distribution of latent variable  $z_{pr}$  by a linear projection.

$$p(z_{pr}|c) = N(z_{pr}; \mu_{pr\_p}(c), \sigma_{pr\_p}(c)) \quad (4)$$

## 2.2 BN Constraint Module

In this module, the BN encoder extracts the latent representation of pronunciation information  $z_{pr}$  from the BN feature and produces the posterior normal distribution  $q(z_{pr}|BN)$  with the mean  $\mu_{pr\_q}$  and variance  $\sigma_{pr\_q}$ . BN feature is usually the feature map of a neural network layer. Specifically, the BN adopted in this paper is the output of an ASR encoder, which is generally considered to contain only linguistic and prosodic information such as pronunciation, intonation, accent, and very limited speaker information [11]. The ASR model is usually

trained with a large multi-speaker multi-condition dataset, and the BN feature extracted by it is also believed to be noise-robust and speaker-independent. The BN encoder consists of multiple layers of Conv1d, ReLU activation, Layer Normalization, Dropout, and a layer of linear projection to produce the mean and variance.

### 2.3 Prior Encoder

The BN decoder in the prior encoder module also acts as the decoder of the first CVAE structure. Given the speaker identity condition  $spk$ , the BN decoder extracts the latent representation of acoustic feature from sampled  $z_{pr}$  and generates the prior normal distribution with mean  $\mu_\theta$  and variance  $\sigma_\theta$  of  $z_{ac}$ . Following VITS, a normalizing flow [4, 17]  $f_\theta$  is added to the prior encoder to improve the expressiveness of the prior distribution of the latent variable  $z_{ac}$ .

$$p(f_\theta(z_{ac})|z_{pr}, spk) = N(f_\theta(z_{ac}); \mu_\theta(z_{pr}, spk), \sigma_\theta(z_{pr}, spk)) \quad (5)$$

$$p(z_{ac}|z_{pr}, spk) = p(f_\theta(z_{ac})|z_{pr}, spk) \left| \det \frac{\partial f_\theta(z_{ac})}{\partial z_{ac}} \right| \quad (6)$$

### 2.4 Posterior Encoder

The posterior encoder module extracts the latent representation  $z_{ac}$  from the waveform  $y$ . The mel spectrum extractor in it is a fixed signal processing layer without updatable weights. The encoder firstly extracts the mel spectrum from the raw waveform through the signal processing layer. Unlike VITS, the posterior encoder takes the mel spectrum as input instead of the linear spectrum. We use multiple layers of Conv1d, ReLU activation, Layer Normalization, and Dropout to extract a sequence of hidden vector and then produces the mean and variance of the posterior distribution  $q(z_{ac}|y)$  by a Conv1d layer. Then we can get the latent  $z_{ac}$  sampled from  $q(z_{ac}|y)$  using the reparametrization trick.

### 2.5 Decoder

The decoder generates audio waveforms from the intermediate representation  $z_{ac}$ . We use HiFi-GAN generator G [9] as the decoder. For more efficient training, we only feed the sliced  $z_{ac}$  instead of the entire length into the decoder to generate the corresponding audio segment. We also use GAN-based [6] training to improve the quality of the synthesized speech. The discriminator D follows HiFiGAN’s Multi-Period Discriminator (MPD) and Multi-Scale Discriminator (MSD) [9]. Specifically, the GAN losses for the generator G and discriminator D are defined as:

$$L_{adv}(G) = E_{(z_{ac})} [(D(G(z_{ac})) - 1)^2] \quad (7)$$

$$L_{adv}(D) = E_{(y, z_{ac})} [(D(y) - 1)^2 + (D(G(z_{ac})))^2] \quad (8)$$

**Table 1.** Experimental results in terms of subjective mean opinion score (MOS) with confidence intervals of 95% and two objective metrics. Note: the results of VITS-DAT are missing as the system cannot converge properly during training.

| Accent   | Model       | Subjective Evaluation |                  |                  | Objective Evaluation       |               |
|----------|-------------|-----------------------|------------------|------------------|----------------------------|---------------|
|          |             | SMOS↑                 | NMOS↑            | AMOS↑            | Speaker cosine Similarity↑ | Duration MAE↓ |
| Shanghai | T2B2M       | 3.85±0.03             | 3.68±0.06        | 3.73±0.05        | 0.76                       | 3.51          |
|          | Accent-VITS | <b>3.86±0.02</b>      | <b>3.76±0.02</b> | <b>3.76±0.06</b> | <b>0.83</b>                | <b>3.06</b>   |
| Henan    | T2B2M       | 3.79±0.07             | 3.75±0.03        | 3.59±0.02        | 0.78                       | 3.65          |
|          | Accent-VITS | <b>3.87±0.02</b>      | <b>3.92±0.04</b> | <b>3.62±0.04</b> | <b>0.81</b>                | <b>3.21</b>   |
| Dongbei  | T2B2M       | 3.88±0.05             | 3.87±0.03        | 3.50±0.03        | 0.79                       | 3.49          |
|          | Accent-VITS | <b>4.01±0.02</b>      | <b>4.14±0.05</b> | <b>3.88±0.02</b> | <b>0.87</b>                | <b>3.01</b>   |
| Sichuan  | T2B2M       | 3.94±0.06             | <b>3.82±0.04</b> | 3.69±0.06        | <b>0.84</b>                | 3.55          |
|          | Accent-VITS | <b>4.06±0.04</b>      | 3.80±0.02        | <b>3.81±0.06</b> | <b>0.84</b>                | <b>3.14</b>   |
| Average  | T2B2M       | 3.87±0.05             | 3.78±0.04        | 3.63±0.02        | 0.79                       | 3.55          |
|          | Accent-VITS | <b>3.95±0.04</b>      | <b>3.91±0.03</b> | <b>3.77±0.06</b> | <b>0.84</b>                | <b>3.11</b>   |

## 2.6 Duration Predictor

In the training process, the LR module expands the phoneme-level representation using the ground truth duration, denoted as  $d$ . In the inference process, the LR module expands the representation using the predicted duration, denoted as  $\hat{d}$ , obtained from the duration predictor. Unlike VITS, we utilize a duration predictor consisting of multiple layers of Conv1d, ReLU activation, Layer Normalization, and Dropout instead of a stochastic duration predictor due to the significant correlation between accent-specific pronunciation prosody and duration information [23]. For the duration loss  $L_{dur}$ , we use MSE loss between  $\hat{d}$  and  $d$ .

## 2.7 Final Loss

With the above hierarchical CVAE and adversarial training, we optimize our proposed model with the full objective:

$$L = L_{adv}(G) + L_{fm}(G) + L_{cvae} + \lambda L_{dur} \quad (9)$$

$$L(D) = L_{adv}(D) \quad (10)$$

where  $L_{adv}(G)$  and  $L_{adv}(D)$  are the GAN loss of G and D respectively, and feature matching loss  $L_{fm}$  is added to improve the stability of the training. The  $L_{cvae}$  consists of the reconstruction loss and two KL losses.

### 3 Experiments

#### 3.1 Datasets

The experimental data consists of high-quality standard Mandarin speech data and accent Mandarin speech data from four different regions: Sichuan, Dongbei (Northeast China), Henan, and Shanghai. Specifically, we use DB1<sup>4</sup> as the high-quality standard Mandarin data which contains 10,000 utterances recorded in a studio from a professional female anchor. The total duration is approximately 10.3 hours. The accent data from the four regions were also recorded in a recording studio by speakers from these regions. Among them, Sichuan, Dongbei, and Shanghai each have two speakers, one male and one female respectively, while Henan has only one female speaker. In detail, Sichuan, Dongbei, Henan, and Shanghai accent data have 2794, 3947, 2049, and 4000 utterances, respectively. The duration of the accent data is approximately 13.7 hours. So we have a total of 4 accents and 8 speakers in our training data.

All the audio recordings are downsampled to 16kHz. We utilize 80-dim mel-spectrograms with 50ms frame length and 12.5ms frame shift. Our ASR model is based on the WeNet U2++ model [21] trained on 10,000 hours of data from the WenetSpeech corpus [22]. We use the Conformer-based encoder output as our BN feature with 512-dim. The BN feature is further interpolated to match the sequence length of the mel-spectrogram.

**Table 2.** Results of ablation studies.

| Accent  | Model         | Subjective Evaluation |                  |                  | Objective Evaluation       |               |
|---------|---------------|-----------------------|------------------|------------------|----------------------------|---------------|
|         |               | SMOS↑                 | NMOS↑            | AMOS↑            | Speaker cosine Similarity↑ | Duration MAE↓ |
| Average | Accent-VITS   | <b>3.95±0.04</b>      | <b>3.91±0.03</b> | <b>3.77±0.06</b> | <b>0.84</b>                | <b>3.11</b>   |
|         | -BN encoder   | 3.82±0.04             | 3.72±0.01        | 3.61±0.02        | 0.77                       | 3.41          |
|         | -BN decoder   | 3.89±0.06             | 3.86±0.03        | 3.66±0.03        | 0.82                       | 3.13          |
|         | -BN (enc,dec) | 3.05±0.07             | 3.24±0.04        | 2.93±0.07        | 0.69                       | 4.03          |

#### 3.2 Model Configuration

We implemented the following three models for comparison.

- Text2BN2Mel (T2B2M) [2, 24]: a three-stage accent transfer system composed of independently trained models for Text2BN, BN2Mel, and neural

<sup>4</sup> [https://www.data-baker.com/open\\_source.html](https://www.data-baker.com/open_source.html)

Vocoder. The Text2BN model predicts BN feature that contains accent pronunciation and content information from the input text. The BN2Mel model predicts mel-spectrogram based on the input BN feature and speaker identity. HiFiGAN V1 is used as the vocoder.

- VITS-DAT: an accent transfer model based on VITS and DAT. For disentangling accent and speaker timbre information, we add a DAT module composed of a gradient reversal layer and a speaker classifier to the output of the text encoder in VITS. Note that we also use a non-stochastic duration predictor and LR module instead of a stochastic duration predictor and MAS [7] in this model.
- Accent-VITS: the proposed accent transfer model in this paper.

In the text frontend processing, we assigned a different phoneme set for each accent or standard Mandarin. The ground truth phoneme-level duration of all datasets is extracted by force-alignment tools. The above comparison models are trained for 400k steps. The batch size of all the models is 24. The initial learning rate of all the models is  $2e-4$ . The Adam optimizer with  $\beta_1 = 0.8$ ,  $\beta_2 = 0.99$  and  $\epsilon = 10^{-9}$  is used to train all them.

### 3.3 Subjective Evaluation

The TTS test set consists of both short and long total of 30 sentences for each accent without overlap with the training set. Each test sentence was synthesized by combining all the speakers in the training data separately. The VITS-DAT system can not converge due to the instability of the combination of variational inference and DAT. Therefore we do not evaluate the results of VITS-DAT here. We randomly selected 20 synthetic utterances for each accent, resulting in a total of 80 utterances for subjective listening. We asked fifteen listeners for each accent to assess speaker similarity (SMOS) and speech naturalness (NMOS). There are twenty local accent listeners to evaluate the accent similarity (AMOS), with five listeners for each accent. Particularly for the SMOS test, we use target speakers' real recordings as reference. The results are summarized in Table 1.

**Speaker Similarity.** The results shown in Table 1 indicate that Accent-VITS can achieve the best performance in speaker similarity. Among the transfer results of all four accents, the SMOS score of Accent-VITS is better than that of T2B2M. This shows that our end-to-end model Accent-VITS effectively avoids the error accumulation and mismatch problems in the multi-stage model so that it can synthesize speech with more realistic target speaker timbre.

**Speech Naturalness.** In the NMOS test, we ask the listeners to pay more attention to the general prosody such as rhythm and expressiveness of the audio. From Table 1 we can see that the NMOS score of Accent-VITS is very close to T2B2M on the Sichuan accent and outperforms T2B2M on the other three accents. This indicates that the speech synthesized by the end-to-end model Accent-VITS is more natural than T2B2M on average.



**Accent Similarity.** In the AMOS test, we ask accent listeners to assess the similarity between synthesized speech and target accent, ignoring the naturalness of general prosody. The results in Table 1 show that Accent-VITS achieves a higher AMOS score than T2B2M, which indicates that Accent-VITS can model accent attribute information better than T2B2M thanks to its hierarchical modeling of accent information and acoustic features.

### 3.4 Objective Evaluation

Objective metrics, including speaker cosine similarity and duration mean absolute error (Duration MAE), are also calculated.

**Speaker Cosine Similarity.** We calculate the cosine similarity on the generated samples to further verify the speaker similarity. Specifically, we train an ECAPA-TDNN model [3] using 6000 hours of speech from 18083 speakers to extract x-vectors. The cosine similarity to the target speaker audio is measured on all synthetic utterances. The results are also shown in Table 1. The speaker cosine similarity score of Accent-VITS is also higher than that of T2B2M on three accents except for the Sichuan accent. Compared with the T2B2M, Accent-VITS gets higher scores of speaker cosine similarity in Shanghai, Henan, and Dongbei accents. And in the Sichuan accent, both are equal. This further demonstrates that Accent-VITS is better than T2B2M in modeling the target speaker timbre.

**Duration MAE.** Prosody variations are key attributes of accent rendering, which is largely reflected in the perceived duration of pronunciation units. Therefore, we further calculate the duration mean absolute error between the predicted duration results of different models and the ground truth. The results in Table 1 show that Accent-VITS gets lower Duration MAE scores than T2B2M in transfer results of all four accents, which means that the transfer results of Accent-VITS are closer to the target accent in prosody than the transfer results of T2B2M.

### 3.5 Ablation Study

To investigate the importance of our proposed methods in Accent-VITS, three ablation systems were obtained by dropping the BN encoder and BN decoder respectively, and dropping both of them simultaneously, referred to as *-BN encoder*, *-BN decoder*, and *-BN (enc, dec)*. When dropping the BN encoder alone, the FFT blocks module directly predicts BN as an intermediate representation. We use MSE loss between the predicted BN and the ground truth BN as the constraint. The BN decoder module takes BN as input. When dropping the BN decoder alone, the distribution of  $z_{pr}$  directly as the prior distribution of  $z_{ac}$ . The flow module takes the sampled  $z_{pr}$  as input. When dropping both of them simultaneously, the FFT blocks module predicts the distribution of BN as the prior distribution of  $z_{ac}$ . We use MSE loss between the sampled BN and the ground truth BN as the constraint.

The results of ablation studies are shown in Table 2. As can be seen, dropping these methods brings performance degradation in terms of subjective evaluation and objective evaluation. Especially dropping both of them simultaneously leads to significantly performance degradation. This validates the effectiveness and importance of the hierarchical CVAE modeling structure in our proposed model.

## 4 CONCLUSIONS

In this paper, we propose Accent-VITS, a VITS-based end-to-end model with a hierarchical CVAE structure for accent transfer. The hierarchical CVAE respectively models accent pronunciation information with the constraint of BN and acoustic features with the constraint of mel-spectrum. Experiments on professional Mandarin data and accent data show that Accent-VITS significantly outperforms the Text2BN2Mel+Neural-Vocoder three-stage approach and the VITS-DAT approach.

## References

1. Acuna, D., Law, M.T., Zhang, G., Fidler, S.: Domain adversarial training: A game perspective. In: The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022. OpenReview.net (2022)
2. Dai, D., Chen, Y., Chen, L., Tu, M., Liu, L., Xia, R., Tian, Q., Wang, Y., Wang, Y.: Cloning one’s voice using very limited data in the wild. In: IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2022, Virtual and Singapore, 23-27 May 2022. pp. 8322–8326. IEEE (2022)
3. Desplanques, B., Thienpondt, J., Demuynck, K.: ECAPA-TDNN: emphasized channel attention, propagation and aggregation in TDNN based speaker verification. In: Meng, H., Xu, B., Zheng, T.F. (eds.) Interspeech 2020, 21st Annual Conference of the International Speech Communication Association, Virtual Event, Shanghai, China, 25-29 October 2020. pp. 3830–3834. ISCA (2020)
4. Dinh, L., Sohl-Dickstein, J., Bengio, S.: Density estimation using real NVP. In: 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings. OpenReview.net (2017)
5. Ganin, Y., Ustinova, E., Ajakan, H., Germain, P., Larochelle, H., Laviolette, F., Marchand, M., Lempitsky, V.S.: Domain-adversarial training of neural networks. *J. Mach. Learn. Res.* **17**, 59:1–59:35 (2016)
6. Goodfellow, I.J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A.C., Bengio, Y.: Generative adversarial nets. In: Ghahramani, Z., Welling, M., Cortes, C., Lawrence, N.D., Weinberger, K.Q. (eds.) *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014*, December 8-13 2014, Montreal, Quebec, Canada. pp. 2672–2680 (2014)
7. Kim, J., Kong, J., Son, J.: Conditional variational autoencoder with adversarial learning for end-to-end text-to-speech. In: Meila, M., Zhang, T. (eds.) *Proceedings of the 38th International Conference on Machine Learning, ICML 2021*, 18-24 July 2021, Virtual Event. *Proceedings of Machine Learning Research*, vol. 139, pp. 5530–5540. PMLR (2021)

8. Kolluru, B., Wan, V., Latorre, J., Yanagisawa, K., Gales, M.J.F.: Generating multiple-accent pronunciations for TTS using joint sequence model interpolation. In: Li, H., Meng, H.M., Ma, B., Chng, E., Xie, L. (eds.) INTERSPEECH 2014, 15th Annual Conference of the International Speech Communication Association, Singapore, September 14-18, 2014. pp. 1273–1277. ISCA (2014)
9. Kong, J., Kim, J., Bae, J.: Hifi-gan: Generative adversarial networks for efficient and high fidelity speech synthesis. In: Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., Lin, H. (eds.) Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual (2020)
10. Lee, S., Kim, S., Lee, J., Song, E., Hwang, M., Lee, S.: Hierspeech: Bridging the gap between text and speech by hierarchical variational inference using self-supervised representations for speech synthesis. In: NeurIPS (2022)
11. Li, J., Deng, L., Gong, Y., Haeb-Umbach, R.: An overview of noise-robust automatic speech recognition. IEEE ACM Trans. Audio Speech Lang. Process. **22**(4), 745–777 (2014). <https://doi.org/10.1109/TASLP.2014.2304637>, <https://doi.org/10.1109/TASLP.2014.2304637>
12. Liu, R., Sisman, B., Gao, G., Li, H.: Controllable accented text-to-speech synthesis. CoRR **abs/2209.10804** (2022). <https://doi.org/10.48550/arXiv.2209.10804>, <https://doi.org/10.48550/arXiv.2209.10804>
13. Liu, S., Yang, S., Su, D., Yu, D.: Referee: Towards reference-free cross-speaker style transfer with low-quality data for expressive speech synthesis. In: IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2022, Virtual and Singapore, 23-27 May 2022. pp. 6307–6311. IEEE (2022)
14. Loots, L., Niesler, T.: Automatic conversion between pronunciations of different english accents. Speech Commun. **53**(1), 75–84 (2011)
15. de Mareüil, P.B., Vieru-Dimulescu, B.: The contribution of prosody to the perception of foreign accent. Phonetica **63**(4), 247–267 (2006)
16. Ren, Y., Ruan, Y., Tan, X., Qin, T., Zhao, S., Zhao, Z., Liu, T.: Fastspeech: Fast, robust and controllable text to speech. In: Wallach, H.M., Larochelle, H., Beygelzimer, A., d’Alché-Buc, F., Fox, E.B., Garnett, R. (eds.) Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada. pp. 3165–3174 (2019)
17. Rezende, D.J., Mohamed, S.: Variational inference with normalizing flows. In: Bach, F.R., Blei, D.M. (eds.) Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015. JMLR Workshop and Conference Proceedings, vol. 37, pp. 1530–1538. JMLR.org (2015)
18. Shu, R., Bui, H.H., Narui, H., Ermon, S.: A DIRT-T approach to unsupervised domain adaptation. In: 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings. OpenReview.net (2018)
19. Sun, L., Li, K., Wang, H., Kang, S., Meng, H.M.: Phonetic posteriorgrams for many-to-one voice conversion without parallel data training. In: IEEE International Conference on Multimedia and Expo, ICME 2016, Seattle, WA, USA, July 11-15, 2016. pp. 1–6. IEEE Computer Society (2016)
20. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. In: Guyon, I., von Luxburg, U., Bengio, S., Wallach, H.M., Fergus, R., Vishwanathan, S.V.N., Garnett, R. (eds.) Advances

- in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA. pp. 5998–6008 (2017)
21. Yao, Z., Wu, D., Wang, X., Zhang, B., Yu, F., Yang, C., Peng, Z., Chen, X., Xie, L., Lei, X.: Wenet: Production oriented streaming and non-streaming end-to-end speech recognition toolkit. In: Hermansky, H., Cernocký, H., Burget, L., Lamel, L., Scharenborg, O., Motlíček, P. (eds.) Interspeech 2021, 22nd Annual Conference of the International Speech Communication Association, Brno, Czechia, 30 August - 3 September 2021. pp. 4054–4058. ISCA (2021)
  22. Zhang, B., Lv, H., Guo, P., Shao, Q., Yang, C., Xie, L., Xu, X., Bu, H., Chen, X., Zeng, C., Wu, D., Peng, Z.: WENETSPEECH: A 10000+ hours multi-domain mandarin corpus for speech recognition. In: IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2022, Virtual and Singapore, 23-27 May 2022. pp. 6182–6186. IEEE (2022)
  23. Zhang, Y., Cong, J., Xue, H., Xie, L., Zhu, P., Bi, M.: Visinger: Variational inference with adversarial learning for end-to-end singing voice synthesis. In: IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2022, Virtual and Singapore, 23-27 May 2022. pp. 7237–7241. IEEE (2022)
  24. Zhang, Y., Wang, Z., Yang, P., Sun, H., Wang, Z., Xie, L.: Accentspeech: Learning accent from crowd-sourced data for target speaker TTS with accents. In: Lee, K.A., Lee, H., Lu, Y., Dong, M. (eds.) 13th International Symposium on Chinese Spoken Language Processing, ISCSLP 2022, Singapore, December 11-14, 2022. pp. 76–80. IEEE (2022)
  25. Zhou, X., Zhang, M., Zhou, Y., Wu, Z., Li, H.: Accented text-to-speech synthesis with limited data. CoRR **abs/2305.04816** (2023). <https://doi.org/10.48550/arXiv.2305.04816>, <https://doi.org/10.48550/arXiv.2305.04816>