



Adversarial Multi-Task Learning for Disentangling Timbre and Pitch in Singing Voice Synthesis

Tae-Woo Kim, Min-Su Kang, Gyeong-Hoon Lee

Speech AI Lab., AI Center, NCSOFT Corp., Republic of Korea

{ktw0114, mskang, ghlee0304}@ncsoft.com

Abstract

Recently, deep learning-based generative models have been introduced to generate singing voices. One approach is to predict the parametric vocoder features consisting of explicit speech parameters. This approach has the advantage that the meaning of each feature is explicitly distinguished. Another approach is to predict mel-spectrograms for a neural vocoder. However, parametric vocoders have limitations of voice quality and the mel-spectrogram features are difficult to model because the timbre and pitch information are entangled. In this study, we propose a singing voice synthesis model with multi-task learning to use both approaches – acoustic features for a parametric vocoder and mel-spectrograms for a neural vocoder. By using the parametric vocoder features as auxiliary features, the proposed model can efficiently disentangle and control the timbre and pitch components of the mel-spectrogram. Moreover, a generative adversarial network framework is applied to improve the quality of singing voices in a multi-singer model. Experimental results demonstrate that our proposed model can generate more natural singing voices than the single-task models, while performing better than the conventional parametric vocoder-based model.

Index Terms: singing voice synthesis, adversarial training, multi-task learning, timbre, pitch

1. Introduction

Singing voice synthesis (SVS) is a generative model to synthesize singing voices according to musical scores and lyrics. Although the musical scores provide note pitch and duration information, it is difficult to model singing voices as they have a longer vowel duration and wider pitch range than speech signals in general. Moreover, singing voices have timbre, which is everything except the pitch or loudness [1], and pitch information that can change according to the singer's vocalizations and expressions. The singer's timbral characteristics depend on the voice source and vocal tract [2]. To model natural singing voices, SVS is modeled considering the interdependence between timbre and pitch.

In recent years, the designs of most SVS models [3–10] have been inspired by the neural network-based architecture in text-to-speech (TTS) models [11–17]. Recent Korean SVS systems [6–8] are able to generate mel-spectrograms directly by disentangling the timbre and pitch related spectra without external signal processing to generate realistic and natural singing voice. Lee et al. [6] have proposed an adversarial trained end-to-end Korean SVS system that is based on Deep Convolutional TTS [17] and includes a phonetic enhancement masking decoder that predicts the formant spectral mask. In [7], the singer's identity is conditioned on each independent decoder of [6] to generate the timbre and singing style. N-Singer [8], a non-autoregressive SVS model, independently models the phoneme

and pitch modules to generate accurately pronounced singing voices. However, in these methods [6–8], linguistic and note information is fed independently in each module to separate the timbre and pitch representations in an unsupervised manner. Additionally, they have limitations in generating natural singing voices because timbre and pitch are assumed to be independent of each other.

In this study, we propose an adversarial multi-task learning-based SVS model to disentangle the timbre and pitch representations. The main task of the proposed method is to predict the mel-spectrogram, while the auxiliary task is to predict the features of timbre and pitch. The proposed model has been trained in two phases and adversarial trained using discriminators in both phases. In the first phase, it is pre-trained on the auxiliary task. This pre-training allows the two decoders to represent timbre and pitch, respectively. In the second phase, it is jointly trained on the main and auxiliary tasks. Thus, the disentangled mel-spectrograms for timbre and pitch are predicted by the respective decoders and integrated into the final mel-spectrogram. In the experiments, the proposed model performs better than the single-task SVS models in terms of perceptual quality and naturalness. Further, synthesizing the timbre and pitch features as predicted by the auxiliary task using the WORLD vocoder [18] performs better than the conventional WORLD vocoder-based SVS model.

2. Related works

Multi-task learning (MTL) has been widely used in speech applications related to enhancement [19, 20], recognition [21] and synthesis [22, 23]. It is a learning paradigm that improves the performance of generalizations by sharing related knowledge from jointly learning multiple tasks [24]. The model proposed in this study applies MTL to the SVS task to jointly learn two tasks – WORLD vocoder and neural vocoder feature prediction.

Generative adversarial networks (GAN) are a powerful method to solve over-smoothing problems and generate high-resolution images [25]. In speech synthesis, Yang et al. [26] improved speech quality by jointly training a generator and a conditional discriminator using the speaker embeddings in multi-speaker models. We adopt the conditional GAN framework using singer embeddings for the multi-singer SVS model, as in [26].

3. Proposed methods

As depicted in Figures 1(a) and (b), the generators of the single-task SVS model and proposed multi-task SVS model consist of two encoders, two decoders, and a mel-predictor. The single-task model is based on the architecture of N-Singer [8] wherein phonemes and note pitches are independently modeled. Unlike the single-task model, the proposed model's encoder outputs are

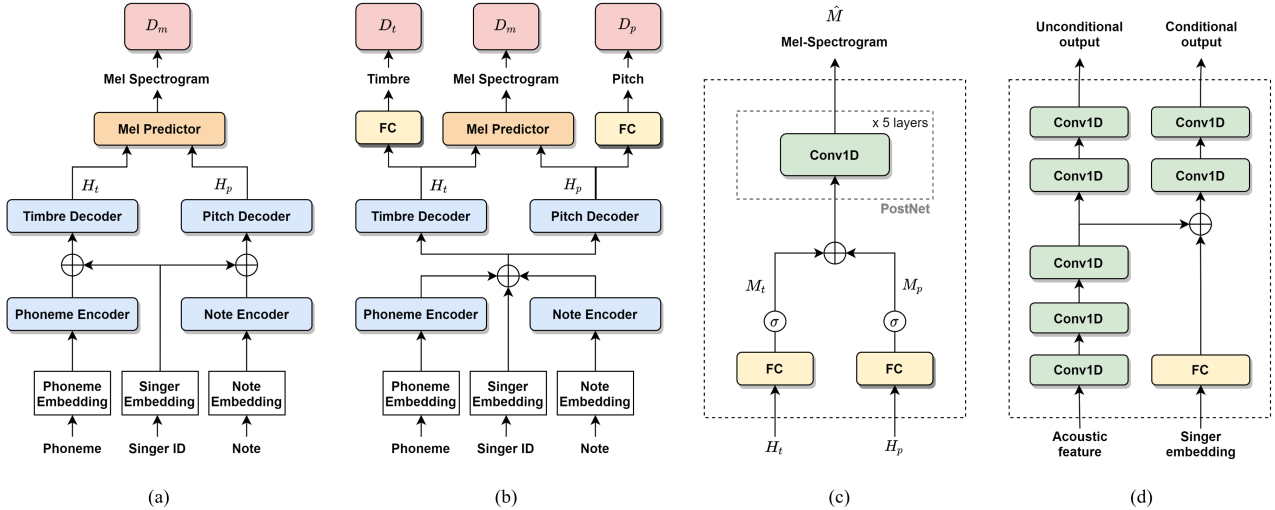


Figure 1: Architecture overview. (a) is a single-task SVS model, (b) is the proposed multi-task SVS model, (c) is a mel-predictor, and (d) is a discriminator.

integrated such that the timbre and pitch decoders are guided through the auxiliary tasks, thereby enabling them to represent each characteristic. Then, predicted timbre and pitch representations are integrated by the mel-predictor to predict the final mel-spectrogram. Thus, the mel-spectrogram and auxiliary features are adversarial trained with discriminators.

3.1. Embeddings

As with previous Korean SVS systems [6–10], we utilize the advantage that one syllable in the lyrics matches one note of musical instrument digital interface (MIDI). First, we decompose Korean syllables into a phoneme sequence using the Korean grapheme-to-phoneme algorithm. Second, we get a frame-level phoneme sequence $T \in \mathbb{R}^{1 \times L}$ by assigning phoneme to frame corresponding to each phoneme using a MIDI note sequence $P \in \mathbb{R}^{1 \times L}$ that includes starting time, duration, and pitch, where L is the length of the acoustic feature. Third, because each Korean syllable consists of an onset, a nucleus, and an optional coda in general, we assign the onset and coda to the first and last three frames, respectively, while assigning the remaining frames to the nucleus within the interval corresponding to each syllable in the phoneme sequence. The input sequences are embedded into D -dimensional dense vectors $E_T \in \mathbb{R}^{D \times L}$ and $E_P \in \mathbb{R}^{D \times L}$ using learnable lookup tables, respectively. The embedded phoneme sequence E_T and note sequence E_P are passed on to the phoneme encoder and note encoder respectively, after which they are integrated with the singer embedding vector $E_S \in \mathbb{R}^{D \times 1}$. The phoneme and note encoders consist of stacked conformer blocks [27]. The singer embedding vector E_S is obtained using a learnable lookup table.

3.2. Timbre and pitch decoders

The two decoders in our model consist of stacked conformer blocks to predict the timbre and pitch representations. As the two decoders cannot disentangle the individual representations of timbre and pitch in an unsupervised manner owing to the common input received from the encoders, an MTL approach is introduced. With the auxiliary task of explicitly predicting tim-

bre and pitch features, the two decoder networks are shared to learn the timbre and pitch representations, respectively. First, the timbre representation H_T , which is the output of the shared timbre decoder, is passed through the fully connected (FC) layer and the mel-predictor. The FC layer predicts the timbre features – mel-generalized cepstrum (MGC), band aperiodicity (BAP), and voiced/unvoiced (V/UV) flags. On the other side, the pitch representation H_P , which is the output of the shared pitch decoder, is passed to the other FC layer predicting log-scale F0 (logF0) and to the mel-predictor. During training, the ground truth V/UV flags are applied to train the logF0 of the voiced section. Thus, the loss for the pitch is calculated as

$$L_{pitch} = L_1(p, \hat{p} \odot v), \quad (1)$$

where \odot , p , \hat{p} , and v are the element-wise multiplication, target logF0, predicted logF0, and ground-truth V/UV flags, respectively. Then, the loss for all the auxiliary features is as follows:

$$L_{aux} = w_m \cdot L_{mgc} + w_b \cdot L_{bap} + w_v \cdot L_{vuv} + w_p \cdot L_{pitch}, \quad (2)$$

where L_{mgc} , L_{bap} and L_{vuv} indicate the loss of the MGC, BAP, and V/UV, respectively. L_{vuv} uses the binary cross-entropy function, while the others use the L1 loss function. w_m , w_b , w_v , and w_p are the respective scalar weights corresponding to each loss.

3.3. Mel-predictor

As shown in Figure 1(c), the timbre and pitch representations that are output from the two decoders are used to predict the log-scale mel-spectrogram by the mel-predictor. First, the two hidden representations pass through the two respective FC layers in Figure 1(c) to reduce dimensionality, after which the sigmoid is applied and output as the visible and interpretable representations, M_T and M_P . In our preliminary experiments, we found that the two representations equal the spectral envelope and pitch harmonics, respectively. According to the source filter theory [28], the pitch harmonics are multiplied by the spec-

tral envelope in the linear spectral domain. However, since this study deals with log-scale mel-spectrograms, multiplication is replaced by a summation operation. Therefore, M_T and M_P are summed and passed through a convolutional neural network-based postnet [16] to predict the final mel-spectrogram. Since the source filter’s modeling is applied to the mel-spectrogram domain, the postnet in the mel-predictor makes the predicted coarse mel-spectrogram close to the true mel-spectrogram M . Consequently, the loss for the mel-spectrograms is represented as:

$$L_{mel} = L_1(M, M_T + M_P) + L_1(M, \hat{M}) \quad (3)$$

3.4. Discriminators

To improve the perceptual quality of the generated singing voices, we adopt the conditional GAN framework proposed in [26]. As shown in Figure 1(d), we use the joint conditional and unconditional discriminator of [26] with the singer embeddings. Unlike previous models [26] that have only used the GAN framework to enhance the mel-spectrogram, we use three discriminators for the mel-spectrogram, MGC and logF0 to further disentangle the timbre and pitch representations. We have experienced that training for the auxiliary features using reconstruction loss is not sufficient to disentangle the timbre and pitch representations. Moreover, since MGC and logF0 in the multi-singer model have high diversities based on prosodic attributes such as the singer’s identity, purely reconstruction loss-based training would weaken the capabilities of MGC and logF0 prediction. That would result in H_t and H_p having entangled attributes between the timbre and pitch due to the powerful mel-predictor, leading to a degradation of performance. Thus, we use the least-squares loss function [29] and the additional feature matching loss L_{fm} for the adversarial training as follows:

$$L_{dis} = \frac{1}{|S_F|} \sum_{i \in S_F} \left[\frac{1}{2} \mathbb{E}_s [D_i(\hat{x}_i)^2 + D_i(\hat{x}_i, s)^2] + \frac{1}{2} \mathbb{E}_{(x,s)} [(D_i(x_i) - 1)^2 + (D_i(x_i, s) - 1)^2] \right], \quad (4)$$

$$L_{adv} = \frac{1}{|S_F|} \sum_{i \in S_F} \mathbb{E}_s [(D_i(\hat{x}) - 1)^2 + (D_i(\hat{x}, s) - 1)^2], \quad (5)$$

$$L_{fm} = \frac{1}{|S_F|} \sum_{i \in S_F} \mathbb{E}_{(x,s)} \left[\sum_{l=1}^L \frac{1}{N_l} \|D_i^{(l)}(x) - D_i^{(l)}(\hat{x})\|_1 \right], \quad (6)$$

where x , \hat{x} , and s are the target feature, generated feature, and singer embedding, respectively. $S_F = \{m, t, p\}$ is the index set including the mel-spectrogram, MGC, and logF0, while L is the total number of layers in each discriminator. Accordingly, the total loss function of the generator as follows:

$$L_{total} = L_{mel} + L_{aux} + L_{adv} + \lambda_{fm} L_{fm} \quad (7)$$

During the first pre-training phase, we remove L_{mel} and set $S_F = \{t, p\}$ to further disentangle the timbre and pitch representations. For the multi-task training phase, we train our model using L_{total} as Equation (7).

4. Experiments

4.1. Experimental setups

4.1.1. Dataset

We collected singing voice data of 50 Korean songs each from both singers (one male and one female) in the children’s song style (CS) and 20 songs each from five female singers in the Korean ballad style (BS). All recordings were sampled at 48kHz with 16-bit quantization. For each singer, the songs were separated into two parts (90% and 10%) for training and testing, respectively. All songs were segmented into a range of 5 to 15 seconds. For the acoustic features, the 80-dimensional log mel-spectrograms were extracted with a fast Fourier transform size, window size, and frame shift of 2,048, 1,920, and 480, respectively. Further, the 60-dimensional MGC, 5-dimensional BAP, F0, and V/UV flag were extracted by WORLD with a hop size of 10 ms. We used additional data from the open source datasets Kiritan [30], CSD [31], and VocalSet [32] to improve the quality of the vocoder.

4.1.2. Model configurations

In our experiments, the phoneme, note, and singer ID were embedded as a 384-dimensional vector. All encoders and decoders were stacked with six conformer blocks [27], each of which were composed of the following modules: multi-head self-attention (MHSA), convolution, and feedforward (two). In the MHSA module, the number of heads and the hidden size of self-attention were 2 and 384, respectively. In the convolution module, the hidden size of pointwise convolutions was 384, and the kernel size of 1D depth wise convolutions was 31. In the feedforward modules, the input/output sizes of the first and second linear layers were 384/1536 and 1536/384, respectively. The dropout rate for each module of the conformer block was 0.1. The FC layers of the mel-predictor described in Figure 1(c) converted the two 384-dimensional hidden representations into timbre and pitch representations of 80 channels. The postnet architecture used is the same as the postnet in [16], where the kernel size and channels are set to 5 and 80, respectively, and all layers except the last layer use the tanh activation function. The discriminators described in Figure 1(d) consist of 1D convolution layers with leaky ReLU activation functions. For the unconditional layers of the discriminators, the number of channels, kernel sizes and strides of 1D convolution layers are [128,256,1024,256,1], [5,9,9,9,5], and [1,2,2,1,1]. For the conditional layers of the discriminators, the parameters of the conditional layers are the same as the 4th and 5th unconditional layers. For the loss weights, we set w_m , w_b , w_v , w_p , and λ_{fm} to 10, 1, 1, 4, and 10, respectively. For the vocoder, we use Parallel WaveGAN (PWG) [33]. To reconstruct the singing voice waveform with a 48k sampling rate, we set the kernel size of each 1D-convolutional layer to 31 to have a large receptive field, while the rest are the same as [5].

4.1.3. Training

In the acoustic model, we first pre-trained for 30,000 iterations to disentangle the representations of timbre and pitch. Then, we trained 270,000 iterations for the multi-task training. In all training phase, the generator is adversarial trained with the discriminators. We used the Adam optimizer ($\beta_1 = 0.9$, $\beta_2 = 0.98$ and $\epsilon = 10^{-9}$) with a batch size of 4. The initial learning rate was set to 0.0001 and halved after every 50,000 iterations. All parameters of the generator and discriminator were initialized

Table 1: The MOS test results with 95% confidence intervals

Method	PWG	WORLD
<i>XiaoiceSing</i>	-	3.35 ± 0.09
<i>N-Singer</i>	3.53 ± 0.09	-
<i>Single-task</i>	3.35 ± 0.08	-
<i>Multi-task (Proposed)</i>	3.76 ± 0.09	3.75 ± 0.08
<i>Ground truth</i>	3.91 ± 0.10	4.23 ± 0.08

via the Xavier initialization [34]. We also trained the vocoder for 1,000,000 iterations using a RAdam [35] optimizer with the same hyperparameters and training method as in [5].

4.2. Results

4.2.1. Listening test

We conducted the mean opinion scores (MOS) listening test on the overall performance including perceptual quality and naturalness. The proposed multi-task SVS model was compared with the single-task SVS model depicted in Figure 1(a) and the N-Singer model. The N-Singer was extended to a multi-singer model by adding singer embeddings to the two encoder outputs. Additionally, the performance of the sample synthesizing the predicted WORLD vocoder features of the proposed model was also compared with *XiaoiceSing*, a conventional WORLD vocoder-based SVS model. Moreover, the *XiaoiceSing* model was extended to a multi-singer model in the same way as above, and was adversarial trained for logF0 and MGC using the discriminators of the proposed model. We included the reconstructed samples (as the ground truth) to check the upper bound of the vocoder’s performance. In each model, 32 samples were prepared by randomly selecting 8 samples each for four singers (two CS singers and two BS singers). Then, we requested 13 native-speaking Korean participants to evaluate the MOS for the overall performance.

In Table1, for the PWG, the MOS results demonstrate that the proposed model performs better than the single-task model and N-Singer. For the WORLD vocoder, the singing voices synthesized from the multi-task model exhibit better performance than the singing voices of the *XiaoiceSing*. It shows that the two decoders are more advantageous than the integrated decoder when modeling timbre and pitch. Although the performance of WORLD vocoder is better than that of PWG for ground truth, the proposed model is better for samples synthesized using the PWG than that of the WORLD. In future work, improving the performance of the 48k sampling rate-based-vocoder for singing voices is left.

4.2.2. Effect of adversarial training for auxiliary task

We conducted an ablation study to confirm the effect of the adversarial training on the auxiliary task. Figure 2 presents an example of a mel-spectrogram comparing two kinds of adversarial training (one with and another without the auxiliary task), where (a) is the synthesized final mel-spectrogram; (b) and (c) are the timbre and pitch representations, respectively. The top line represents adversarial training applied to the mel-spectrogram alone, while the bottom is the adversarially trained mel-spectrogram with MGC and logF0. The bottom panels of (b) and (c) show that the spectral envelope and pitch harmonics are separated. In contrast, the top panels of (b) and (c) show that the pitch harmonic component remains in the timbre rep-

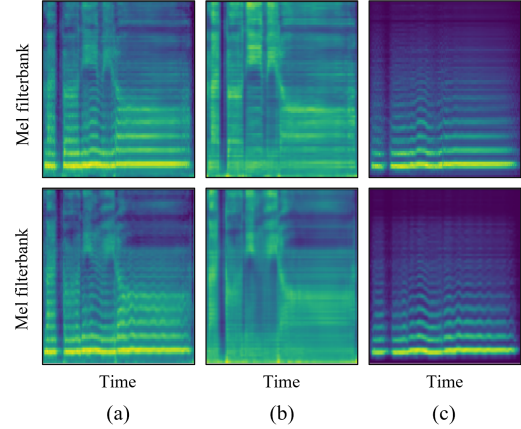


Figure 2: Generated mel-spectrograms of (a) final predicted output, (b) timbre representation, and (c) pitch representation. The top-line figures represent spectrograms from the model trained using adversarial loss for mel-spectrogram alone, while the bottom is from the model using adversarial loss for all features.

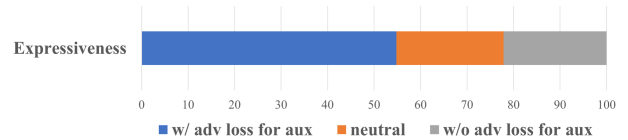


Figure 3: A/B preference test for expressiveness of the generated samples from the proposed model with or without adversarial training for auxiliary task.

resentations, while the pitch representation has a weak vibrato compared to the bottom. Therefore, it shows that the precise predictions of MGC and logF0, which are the auxiliary tasks, help to disentangle the timbre and pitch representations.

We also conducted the A/B preference test to compare the performance difference according to the disentangling ability. For the A/B test, 16 samples of singing voices generated for two BS singers were prepared. Then, we asked the participants to evaluate the expressiveness of the pitch curve and energy dynamics. As shown in Figure 3, 54.8% preferred the proposed model using adversarial loss for the auxiliary task, while 22.1% did not. The remaining 23.1% were neutral. A chi-squared test on the result classes had a p-value of 1.46×10^{-11} , implying that the A/B test result was significant. Therefore, these results confirm that the accurate prediction of MGC and logF0 through the auxiliary tasks helps the decoder better represent the timbre and pitch of singing voices.

5. Conclusions

In this study, we propose an adversarial MTL-based SVS to disentangle the timbre and pitch representations. The proposed approach demonstrates the ability to disentangle timbre and pitch by considering the interdependencies between them. Experimental results confirm that the proposed model performs better than single-task SVS models. Further, the samples synthesized by the auxiliary features also exhibited better performance than that of the conventional WORLD vocoder-based SVS. Some audio samples are available online¹.

¹<https://nc-ai.github.io/speech/publications/amtl-svs/>

6. References

- [1] E. J. Humphrey, S. Reddy, P. Seetharaman, A. Kumar, R. M. Bitner, A. Demetriou, S. Gulati, A. Jansson, T. Jehan, B. Lehner *et al.*, “An introduction to signal processing for singing-voice analysis,” *IEEE Signal ProcESSIng MagazInE*, vol. 1053, no. 5888/19, 2019.
- [2] J. Sundberg and T. D. Rossing, “The science of singing voice,” 1990.
- [3] P. Lu, J. Wu, J. Luan, X. Tan, and L. Zhou, “Xiaoicesing: A high-quality and integrated singing voice synthesis system,” in *Proceedings INTERSPEECH 2020 International Speech Communication Association*, 2020, pp. 1306–1310.
- [4] J. Wu and J. Luan, “Adversarially trained multi-singer sequence-to-sequence singing synthesizer,” in *Proceedings INTERSPEECH 2020 International Speech Communication Association*, 2020, pp. 1296–1300.
- [5] J. Chen, X. Tan, J. Luan, T. Qin, and T.-Y. Liu, “Hifisinger: Towards high-fidelity neural singing voice synthesis,” *arXiv preprint arXiv:2009.01776*, 2020.
- [6] J. Lee, H.-S. Choi, C.-B. Jeon, J. Koo, and K. Lee, “Adversarially trained end-to-end korean singing voice synthesis system,” in *Proc. Interspeech 2019*, 2019, pp. 2588–2592.
- [7] J. Lee, H.-S. Choi, J. Koo, and K. Lee, “Disentangling timbre and singing style with multi-singer singing synthesis system,” in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 7224–7228.
- [8] G.-H. Lee, T.-W. Kim, H. Bae, M.-J. Lee, Y.-I. Kim, and H.-Y. Cho, “N-singer: A non-autoregressive korean singing voice synthesis system for pronunciation enhancement,” in *Proc. Interspeech 2021*, 2021, pp. 1589–1593.
- [9] S. Choi, W. Kim, S. Park, S. Yong, and J. Nam, “Korean singing voice synthesis based on auto-regressive boundary equilibrium gan,” in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 7234–7238.
- [10] J. Tae, H. Kim, and Y. Lee, “Mlp singer: Towards rapid parallel korean singing voice synthesis,” in *2021 IEEE 31st International Workshop on Machine Learning for Signal Processing (MLSP)*. IEEE, 2021, pp. 1–6.
- [11] Y. Ren, Y. Ruan, X. Tan, T. Qin, S. Zhao, Z. Zhao, and T.-Y. Liu, “Fastspeech: Fast, robust and controllable text to speech,” *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [12] Y. Ren, C. Hu, X. Tan, T. Qin, S. Zhao, Z. Zhao, and T.-Y. Liu, “Fastspeech 2: Fast and high-quality end-to-end text to speech,” in *International Conference on Learning Representations*, 2020.
- [13] A. Łańcucki, “Fastpitch: Parallel text-to-speech with pitch prediction,” in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 6588–6592.
- [14] T. Bak, J.-S. Bae, H. Bae, Y.-I. Kim, and H.-Y. Cho, “Fastpitchformant: Source-filter based decomposed modeling for speech synthesis,” *arXiv preprint arXiv:2106.15123*, 2021.
- [15] Y. Wang, R. Skerry-Ryan, D. Stanton, Y. Wu, R. J. Weiss, N. Jaitly, Z. Yang, Y. Xiao, Z. Chen, S. Bengio *et al.*, “Tacotron: Towards end-to-end speech synthesis,” *Proc. Interspeech 2017*, pp. 4006–4010, 2017.
- [16] J. Shen, R. Pang, R. J. Weiss, M. Schuster, N. Jaitly, Z. Yang, Z. Chen, Y. Zhang, Y. Wang, and R. S.-R. *et al.*, “Natural tts synthesis by conditioning wavenet on mel spectrogram predictions,” in *Proceedings ICASSP 2018 IEEE International Conference on Acoustics, Speech and Signal Processing*, 2018, pp. 4779–4783.
- [17] H. Tachibana, K. Uenoyama, and S. Aihara, “Efficiently trainable text-to-speech system based on deep convolutional networks with guided attention,” in *Proceedings ICASSP 2018 IEEE International Conference on Acoustics, Speech and Signal Processing*, 2018, pp. 4784–4788.
- [18] M. Morise, F. Yokomori, and K. Ozawa, “World: a vocoder-based high-quality speech synthesis system for real-time applications,” *IEICE TRANSACTIONS on Information and Systems*, vol. 99, no. 7, pp. 1877–1884, 2016.
- [19] Y. Koizumi, K. Yatabe, M. Delcroix, Y. Masuyama, and D. Takeuchi, “Speech enhancement using self-adaptation and multi-head self-attention,” in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 181–185.
- [20] S. E. Eskimez, X. Wang, M. Tang, H. Yang, Z. Zhu, Z. Chen, H. Wang, and T. Yoshioka, “Human listening and live captioning: Multi-task training for speech enhancement,” *arXiv preprint arXiv:2106.02896*, 2021.
- [21] G. Pironkov, S. Dupont, and T. Dutoit, “Multi-task learning for speech recognition: an overview,” in *ESANN*, 2016.
- [22] Z. Wu, C. Valentini-Botinhao, O. Watts, and S. King, “Deep neural networks employing multi-task learning and stacked bottleneck features for speech synthesis,” in *2015 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2015, pp. 4460–4464.
- [23] R. Liu, B. Sisman, F. Bao, G. Gao, and H. Li, “Modeling prosodic phrasing with multi-task learning in tacotron-based tts,” *IEEE Signal Processing Letters*, vol. 27, pp. 1470–1474, 2020.
- [24] Y. Zhang and Q. Yang, “A survey on multi-task learning,” *IEEE Transactions on Knowledge and Data Engineering*, 2021.
- [25] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial nets,” *Advances in neural information processing systems*, vol. 27, 2014.
- [26] J. Yang, J.-S. Bae, T. Bak, Y. Kim, and H.-Y. Cho, “Ganspeech: Adversarial training for high-fidelity multi-speaker speech synthesis,” *arXiv preprint arXiv:2106.15153*, 2021.
- [27] A. Gulati, J. Qin, C.-C. Chiu, N. Parmar, Y. Zhang, J. Yu, W. Han, S. Wang, Z. Zhang, Y. Wu *et al.*, “Conformer: Convolution-augmented transformer for speech recognition,” in *Proceedings INTERSPEECH 2020 International Speech Communication Association*, 2020, pp. 5036–5040.
- [28] G. Fant, *Acoustic theory of speech production*. Walter de Gruyter, 1970, no. 2.
- [29] X. Mao, Q. Li, H. Xie, R. Y. K. Lau, Z. Wang, and S. P. Smolley, “Least squares generative adversarial networks,” *2017 IEEE International Conference on Computer Vision (ICCV)*, pp. 2813–2821, 2017.
- [30] I. Ogawa and M. Morise, “Tohoku kiritan singing database: A singing database for statistical parametric singing synthesis using japanese pop songs,” *Acoustical Science and Technology*, vol. 42, no. 3, pp. 140–145, 2021.
- [31] S. Choi, W. Kim, S. Park, S. Yong, and J. Nam, “Children’s song dataset for singing voice research,” in *International Society for Music Information Retrieval Conference (ISMIR)*, 2020.
- [32] J. Wilkins, P. Seetharaman, A. Wahl, and B. Pardo, “Vocalset: A singing voice dataset,” in *ISMIR*, 2018, pp. 468–474.
- [33] R. Yamamoto, E. Song, and J.-M. Kim, “Parallel wavegan: A fast waveform generation model based on generative adversarial networks with multi-resolution spectrogram,” in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 6199–6203.
- [34] X. Glorot and Y. Bengio, “Understanding the difficulty of training deep feedforward neural networks,” in *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 2010, pp. 249–256.
- [35] L. Liu, H. Jiang, P. He, W. Chen, X. Liu, J. Gao, and J. Han, “On the variance of the adaptive learning rate and beyond,” *arXiv preprint arXiv:1908.03265*, 2019.