# Adversarially Trained End-to-end Korean Singing Voice Synthesis System

*Juheon Lee[1,2], Hyeong-Seok Choi[1], Chang-Bin Jeon[1,2], Junghyun Koo[1], Kyogu Lee[1,2]*

[1]Music & Audio Research Group, Seoul National University
[2]Center for Super Intelligence, Seoul National University

{juheon2, kekepa15, vinyne}@snu.ac.kr, tonykoo@gmail.com, kglee@snu.ac.kr

## Abstract

In this paper, we propose an end-to-end Korean singing voice synthesis system from lyrics and a symbolic melody using the following three novel approaches: 1) phonetic enhancement masking, 2) local conditioning of text and pitch to the super-resolution network, and 3) conditional adversarial training. The proposed system consists of two main modules; a mel-synthesis network that generates a mel-spectrogram from the given input information, and a super-resolution network that upsamples the generated mel-spectrogram into a linear-spectrogram. In the mel-synthesis network, phonetic enhancement masking is applied to generate implicit formant masks solely from the input text, which enables a more accurate phonetic control of singing voice. In addition, we show that two other proposed methods - local conditioning of text and pitch, and conditional adversarial training - are crucial for a realistic generation of the human singing voice in the super-resolution process. Finally, both quantitative and qualitative evaluations are conducted, confirming the validity of all proposed methods.

**Index Terms**: singing voice synthesis, end-to-end network, phonetic enhancement, conditional adversarial training

## 1. Introduction

With the recent development of deep learning, a learning-based singing voice synthesis (SVS) system, which synthesizes sounds as natural as the concatenative method [1, 2, 3], but can expand more flexibly, is proposed. For example, three SVS systems based on DNN, LSTM, and Wavenet architecture were proposed, respectively [4, 5, 6]. These systems all include an acoustic model that is trained by singing, lyrics, and sheet music paired data, and each acoustical model is trained to predict the vocoder feature used as an input to the vocoder.

Although these neural network-based SVS system can achieve adequate performance, networks predicting vocoder features have limits that cannot exceed the upper bound of vocoder performance. Therefore, it is meaningful to propose an end-to-end framework that directly generates a linear-spectrogram, not a vocoder feature. However, the extension to the end-to-end framework of the SVS system is a challenging task because it involves increased complexity of the model. Creating a more complex target, linear-spectrogram, increases the complexity of the model and requires as much training data to generalize and train these models sufficiently. However, gathering singing audio with aligned lyrics in a controlled environment is a task that requires a lot of effort.

We proposed in this paper a Korean SVS system that can be trained by an end-to-end manner with moderate amounts of data [1]. Our baseline network is designed with the inspiration of DCTTS [7], known as efficiently trainable text to speech

---

[1]The generated result can be found at: ksinging.strikingly.com.

(TTS) system. We applied the following novel approaches to enable end-to-end network training. First, we used the phonetic enhancement masking method, which separately modeled low-level acoustic features related to pronunciation from text information, to make more efficient use of the information contained in the training data. Second, we also proposed a method of reusing input data at the super-resolution stage and training with an adversarial manner to produce better sound quality singing.

The contribution of this paper is as follows: **1)** We designed the end-to-end Korean SVS system and suggested a way to train it effectively. **2)** We proposed a phonetic enhancement masking method that helps to produce more accurate pronunciation. **3)** We proposed a conditional adversarial training method for the generation of more realistic singing voices.

## 2. Related work

The SVS system is similar to the TTS system in terms of synthesizing natural human speech. Recently, the end-to-end TTS system, which is trained as an autoregressive manner, such as Tacotron[8], Deep voice[9], is showing better performance than the conventional method. In addition, various follow-up studies are being conducted that have further controllable elements such as prosody, style, etc [10, 11], or models that can be trained more efficiently [7, 12]. We conducted the study by modifyng the TTS model to suit the SVS task, based on DCTTS[7], which is known to be capable of efficient end-to-end training.

The generative adversarial networks (GAN) is a widely used technique that helps train an arbitrary function to generate a similar sample as the sample from desired data distribution. This training method has been widely accepted in computer vision community and becomes one of the key components to attain photo-realism in super-resolution task. Unlike the success of adversarial training method in image domain, however, only a few works have achieved a reasonable success of training super-resolution task (specifically, band-width extension task) in audio signal processing community [13]. To further leverage the promise of adversarial training in audio generation process, we adopted a few recent works that stabilizes the adversarial training, namely, conditional GAN with projection discriminator [14] and R1 regularization [15] which allow us to jointly train the autoregressive network (mel-synthesis) and super-resolution network making the proposed system as an end-to-end framework.

## 3. Proposed Network

As illustrated in Figure 1, our proposed model consists of two main modules, a mel-synthesis network and a super-resolution network. The mel-synthesis network is trained to produce a mel-spectrogram $M_{1:L}$ from previous mel input $M_{0:L-1}$, time-aligned text $T_{1:L}$, and pitch inputs $P_{1:L}$. With text and
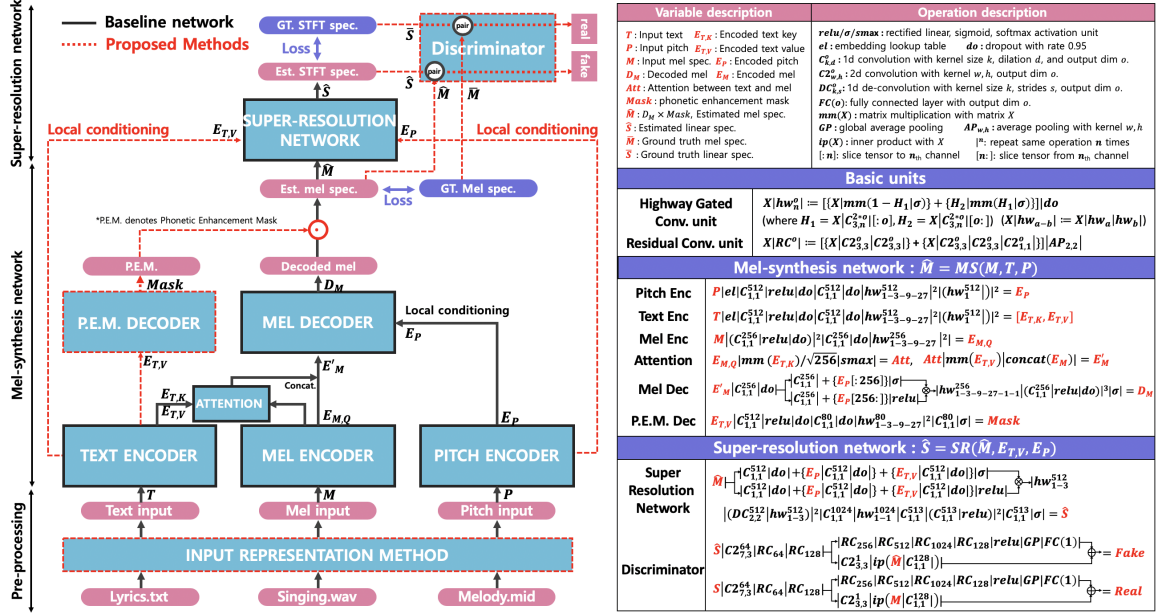
Figure 1: *Proposed system overview (left), detailed structure of each sub-module (right). $X|F|$ denotes $F(X)$.*

pitch information as conditional input, the super-resolution network upsamples the generated mel-spectrogram $M$ to a linear-spectrogram $S$. Finally, the discriminator takes the upsampled result with generated mel-spectrogram to train the network in an adversarial manner. During the test phase, a sequence of mel-spectrogram frames is generated in an autoregressive manner from a given text and pitch input which is then upsampled to linear-spectrogram by super resolution network. Finally, the generated linear-spectrogram is converted to a waveform using Griffin-Lim algorithm [16].
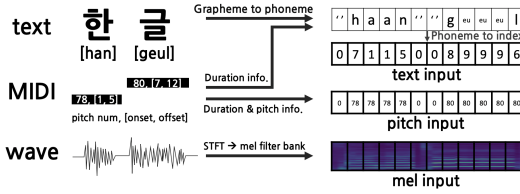


Figure 2: *Input representation method overview*

### 3.1. Input representation

Our training data includes recorded singing voice along with the corresponding text and midi. A single midi note represents pitch information with onset and offset. For the single midi note, one syllable and its corresponding vocal audio section are manually aligned. Figure 2 shows our input representation more concretely. To determine the text input sequence $T \in \mathbb{R}^{1 \times L}$ with length $L$, we referred to the pronunciation system of Korean. A Korean syllable can be decomposed into three phonemes each of which corresponds to onset, nucleus, and coda, respectively. Since the nucleus occupies most of the pronunciation singing in Korean, we assigned onset and coda to the first and the last frame of input text array, respectively, and the rest of the frames with nucleus. Although this does not reflect accurate timing for each phoneme, we empirically found out that a convolution-based network with wide enough range of receptive field can

handle this problem. For pitch input $P \in \mathbb{R}^{1 \times L}$, we simply assigned a pitch number to each frame. In the case of the mel input $M \in \mathbb{R}^{F \times L}$, we used the mel-spectrogram itself, which was extracted from the recorded audio, where $F$ denotes the number of frequency bins.

### 3.2. Mel-synthesis network

The mel-synthesis network $MS(\cdot)$ aims to generate the mel-spectrogram of the next time step from the given text, pitch, and mel input. Based on the text-to-mel network proposed by [7] we modified it to fit the SVS system.

First, in order to enter pitch information, we added pitch encoders with the same structure as text encoders. In addition, the local conditioning method proposed by [17] was used to conduct a conditioning of the encoded pitch on the mel decoder.

Second, we assumed that among the various elements forming a singing voice, information about pronunciation would be able to be controlled independently from text information. We also assumed that if the low-level audio feature that constitutes pronunciation information can be modeled independently, it is possible to focus on the pronunciation information in the data composed of various combinations of pronunciation-pitch, so that training data can be utilized more efficiently to generate more accurate pronounced singing voice. To this end, we designed an additional phonetic enhancement mask decoder, which receives encoded text only as input, and the output of the decoder element-wise multiplied by the output of the mel decoder to create the final mel-spectrogram. As a result, $MS(\cdot)$ can be formulated as follows:

$$\hat{M} = Mask \odot D_M = MS(M, T, P) \qquad (1)$$

We trained the $MS(\cdot)$ network with $L_1$ and binary divergence loss $L_d$ between ground truth and generated mel-spectrogram, and guided attention loss $L_{att}$ as the objective function. Please see [7] for more detailed explanation on the loss terms. We also assumed that the $L_1$ loss between the differential spectrogram $M' = M_{1:L} - M_{0:L-1}$ would be also

beneficial for network to learn more about the relatively short pronounced onset, coda. Therefore, the overall objective function for $MS(\cdot)$ is as follows:

$$L_{MS} = L_1(\hat{M}, M) + L_d(\hat{M}, M) + L_{att} + L_1(\hat{M}', M') \quad (2)$$

### 3.3. Super-resolution network

In this section, we describe the details of the training method for super-resolution network $SR(\cdot)$. The purpose of the SR step is to upsample the generated mel-spectrogram $\hat{M} \in \mathbb{R}^{F \times L}$ into a linear-spectrogram $\hat{S} \in \mathbb{R}^{F' \times L'}$ thereby making it to an audible form, where $F'$ and $L'$ denote the number of frequency bins and temporal bins for linear-spectrogram. The idea of the SR network was proposed in a few previous TTS literatures, including Tacotron and its variants [7, 8]. The major difference between the previous works and our work is twofold. First, we additionally reuse the aligned text and pitch information into the SR network exploiting the useful information in the generation process again. Second, we utilize adversarial training methods to make the SR network produce more realistic sound.

#### 3.3.1. Local conditioning of text and pitch information

Unlike the attention-mechanism based TTS literature, SVS system requires the aligned text and pitch information as inputs for the controllability in the generation process. These information, therefore, can be easily reused in the SR step in the absence of time-alignment process as follows.

$$\hat{S} = SR(\hat{M}, E_{T,V}, E_P) = SR(MS(\cdot), E_{T,V}, E_P) \quad (3)$$

More specifically, each of the output from the text encoder and pitch encoder ($E_{T,V}$ and $E_P$) is fed into a sequence of $1 \times 1$ convolutional and dropout layer [18] which is then fed into a highway network as a local conditioning method as proposed in [17]. For the upsampled $\hat{S}$, SR is trained with the objective function $L_{SR} = L_1(\hat{S}, S) + L_d(\hat{S}, S)$. For the exact network configuration, please refer to Figure 1.

#### 3.3.2. Adversarial training method

Expecting to generate a realistic sound, we adopted a conditional adversarial training method which helps the output distribution of $\hat{S} = SR(\hat{M}, \cdot)$ be similar to the real data distribution $S \sim p(S|M)$. Intuitively, in the conditional adversarial training framework, discriminator $D_\psi$ not only tries to check if $S$ is realistic but also the paired correspondence between $S$ and $M$. Note that, we make a minor assumption that the distribution of $\hat{M} = MS(\cdot)$ approximately follows that of $M$, that is, $p(M) \simeq p(\hat{M})$, allowing the joint training of two modules $MS(\cdot)$ and $SR(\cdot)$. The conditioning to discriminator was done by following [14] with a minor modification. First, the condition $M$ is fed into a 1d-convolutional layer and the intermediate output of discriminator is fed into a $3 \times 3$ 2d-convolutional layer. Then, inner product between the two outputs is done as a projection. Finally, the obtained scalar value is added to the last layer of $D_\psi$ resulting in final logit value. For the exact network configuration please refer to Figure 1.

For the stable adversarial training, a regularization technique on $D_\psi$ has been proposed by several GAN related works [19, 20, 21, 15]. We adopted a simple, yet, effective gradient penalty technique called R1 regularization. This technique penalizes the squared 2-norm of the gradients of $D_\psi$ only when

the sample from true distribution is taken as follows

$$R_1(\psi) = \frac{\gamma}{2} \mathbb{E}_{p(M,S)}[\|\nabla D_\psi(M, S)\|^2]. \quad (4)$$

Note that the output of $D_\psi$ denotes the logit value before the sigmoid function. The final adversarial loss terms ($L_{adv_D}$ and $L_{adv_G}$) for $D_\psi$ and $G_\theta$ are as follows,

$$L_{adv_D}(\theta, \psi) = -\mathbb{E}_{p(M)}[\mathbb{E}_{p(S|M)}[f(D_\psi(M, S))]]$$
$$\quad - \mathbb{E}_{p(\hat{M})}[f(D_\psi((\hat{M}, \hat{S}))] + R_1, \quad (5)$$
$$L_{adv_G}(\theta, \psi) = \mathbb{E}_{p(\hat{M})}[f(D_\psi(\hat{M}, \hat{S}))],$$

where $\theta$ includes not only the parameters of $SR$ but also that of $MS$, hence the two consecutive modules acting as one generator function $G_\theta = SR(MS(\cdot), \cdot)$. The function $f$ is chosen as follows $f(t) = -log(1 + exp(-t))$ resulting in the vanilla GAN loss as in the original GAN paper [22].

## 4. Experiments

### 4.1. Dataset

Since there is no publicly available Korean singing voice dataset, we created the dataset as follows. First, we prepared accompaniment and singing voice MIDI files of 60 Korean pop songs. Next, a professional female vocalist was told to sing to the accompaniment. Then, the singing voice MIDI files were manually realigned so that the recorded audio have the exact alignment with the singing voice MIDI files. Finally, we manually assigned the syllables in lyrics to each MIDI note of singing voice MIDI file. The audio length of the entire dataset excluding the silence is about 2 hours. We used 49 songs for training dataset, 1 song for validation, and 10 songs for test dataset.

### 4.2. Training

We trained the discriminator to minimize $L_{adv_D}$ and the rest of the network to minimize $L_{adv_G}$, $L_{MS}$ and $L_{SR}$. For SR networks, we have to start training after the appropriate level of mel is generated, so we have separately controlled $lr_{SR}$ and $lr_{GAN}$ to add to the objective function. At this point, it was set to $lr_{SR} = \min(0.2 * (\text{iter}/100), 1)$, $lr_{GAN} = \min(0.01 * (\text{int})(\text{iter}/5000), 1)$, respectively.

$$L_{MS,SR} = L_{MS} + lr_{SR} \cdot L_{SR} + lr_{GAN} \cdot L_{adv_G}$$
$$L_D = lr_{GAN} \cdot L_{adv_D} \quad (6)$$

In both cases, we used Adam optimizer [23], which was set to $\beta_1 = 0.5$ and $\beta_2 = 0.9$. The learning rate was scheduled to start from 0.0002 and was halved for every 30,000 iteration. All parameters of the networks were initialized with the Xavier initializer [24].

For the ground truth mel/linear-spectrogram, we first extracted the linear-spectrogram $S$ from audio with $sr = 22050, n_{fft} = 1024, hop = 256$. We then normalized the linear-spectrogram as follows $S \leftarrow (|S|/max(|S|))^\delta$, where $\delta$ denotes a pre-emphasis factor with the value of 0.6 in our case. [2] Afterwards, the mel-spectrogram was obtained by multiplying 80-d of mel filter bank to $S$, and the same normalization method as in $S$ was used. In order to reduce the complexity of the model, we downsampled the mel-spectrogram to the quarter by taking the first frame of every four-frame of the mel-spectrogram giving the relationship of $L' = 4L$.

---

[2]Note that we post emphasized $\hat{S} \leftarrow \hat{S}^{\zeta/\delta}$ where $\zeta$ denotes a post-emphasis factor with the value of 1.3.

### 4.3. Evaluation

We trained a total of five models to see how the three proposed methods - method 1; phonetic enhancement masking, method 2; local conditioning pitch and text to $SR(\cdot)$, method 3; adversarial training method - actually affect the network. The differences between the five models are described in Table 1. 20 audio samples from each model were generated from the test dataset. Apart from the generated samples, we also compared the ground truth samples. **Ground** denotes the actual recorded audio, and **Recons** denotes the reconstructed audio from ground truth magnitude only linear-spectrogram using Griffin-Lim algorithm. Noe that **Recons** samples were included to evaluate the sound quality from the loss of phase information.

Table 1: *Model description and evaluation result.*

| Model | model1 | model2 | model3 | model4 | model5 |
|---|---|---|---|---|---|
| Method | baseline | +(method 1) | +(method 2) | +(method 1,2) | +(method 1,2,3) |

| | Quantitative | | | Qualititative | | |
|---|---|---|---|---|---|---|
| Model | Precision | Recall | F1-score | Pronun.acc | Sound.quality | Naturalness |
| model1 | 0.771 | 0.832 | 0.800 | $2.29 \pm 1.15$ | $2.32 \pm 0.89$ | $2.11 \pm 0.98$ |
| model2 | 0.780 | **0.843** | 0.810 | $2.62 \pm 1.00$ | $2.28 \pm 0.84$ | $2.22 \pm 0.91$ |
| model3 | 0.755 | 0.814 | 0.783 | $2.69 \pm 1.06$ | $2.37 \pm 0.86$ | $2.22 \pm 0.93$ |
| model4 | 0.792 | 0.832 | 0.811 | $2.92 \pm 1.08$ | $2.43 \pm 0.86$ | $2.36 \pm 0.94$ |
| model5 | **0.872** | 0.821 | **0.846** | **$3.23 \pm 1.19$** | **$3.37 \pm 0.94$** | **$3.07 \pm 1.10$** |
| Recons | 0.805 | 0.830 | 0.782 | $4.85 \pm 0.47$ | $4.46 \pm 0.77$ | $4.72 \pm 0.62$ |
| Ground | 0.826 | 0.772 | 0.798 | $4.90 \pm 0.36$ | $4.74 \pm 0.57$ | $4.85 \pm 0.43$ |

#### 4.3.1. Quantitative evaluation

We evaluated whether the network was actually producing a conditioned singing voice for a given input. To do this, we extracted f0 sequence from the generated audio through the world vocoder[25], converted it into a pitch sequence, and compared it to the input pitch sequence. We can judge that the higher the similarity between the two sequence, the more the network generates a singing that reflects the input condition. We calculated the precision, recall and f-score of the generated pitch sequence by frame-wise, and the results are shown in Table 1.

Even in the case of a real recording sample recorded by listening to the original midi accompaniment, it is not easy to adjust the timing and pitch of the correct note, so that a 100% accurate f-score can not be obtained. For all samples that were generated, a f-score similar to or higher than the real recording sample was obtained. This means that the model has generated a singing voice with the correct pitch and timing for at least the real recording for the given input.

#### 4.3.2. Qualitative evaluation

We conducted a listening test to evaluate the quality of the generated singing voice. 19 native Korean speakers were asked to listen to the 20 audio samples from each model. Each participant was asked to evaluate the pronunciation accuracy, sound quality, and naturalness. During the listening test, lyrics of audio samples were provided for more accurate evaluation of pronunciation accuracy. The MOS results are shown in Table 1.

We conducted a paired t-test for each model response and based on this we verified the effectiveness of the proposed methods. For the accuracy of the pronunciation, we obtained significant differences for all comparisons except for models 2 and 3. In other words, all of the proposed methods helped to create more accurate pronunciation singing voices, and the performance was improved to the greatest extent with all three methods. In the case of sound quality, methods 1 and 2 did not significantly affect the improvement, but the applying method 3
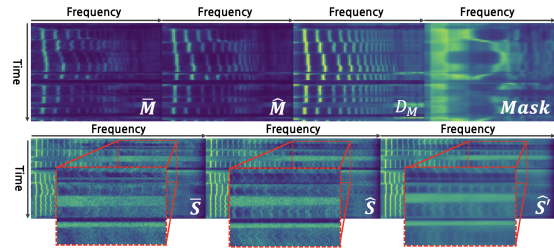


Figure 3: *Generated spectrograms.*

showed a significant increase in score. From this we can confirm that training the network in an adversarial manner improves the quality of the generated audio. Finally, for naturalness, there was a significant improvement when all methods were applied.

### 4.4. Analysis on generated spectrogram

In this section we analyze the features generated by the mel-synthesis and super-resolution networks. In the case of mel-synthesis network, from observing internally generated features, we found that the low-level acoustic feature of pronunciation and pitch could be divided independently without any supervision. From Figure 3, $D_M$ shows the underlying structure of the spectrogram, such as the harmonic structure and the location of f0. In $Mask$, on the other hand, we can observe the shape of determining the intensity of the frequency at every time-step, similar to the feature of the spectral envelope, which contains non-periodic information. This suggests that, from the perspective of source-filter models, one of the techniques that classical speech modelling techniques, our network can generate sources ($D_M$) and filters ($Mask$) separately from frequency domain without any supervised training.

We also analyzed the effect of adversarial training method by observing the generated linear-spectrogram. Three different spectrograms from model4 ($\hat{S}'$: w/o adversarial loss), model5 ($\hat{S}$: w/ adversarial loss), and ground truth spectrogram ($\bar{S}$) are demonstrated in the second row of Figure 3. While $\hat{S}'$ showing the blurry high frequency areas, $\hat{S}$ clearly shows that adversarial training allows the proposed network to generate sample that is closer to the ground truth sample $\bar{S}$. Note that we have confirmed in 4.3.2, listening test that the sound quality can be significantly improved by comparing model4 and model5, which again reinforces our observation.

## 5. Conclusions

In this paper, we proposed the end-to-end Korean singing vocie synthesis system. We showed that using text information to model the phonetic enhancement mask actually worked, and produced more accurate pronunciation. Also, we successfully applied the conditional adversarial training method to the super-resolution stage, which resulted in a higher quality voice.

## 6. Acknowledgements

# 7. References

[1] M. Macon, L. Jensen-Link, E. B. George, J. Oliverio, and M. Clements, "Concatenation-based midi-to-singing voice synthesis," in *Audio Engineering Society Convention 103*. Audio Engineering Society, 1997.

[2] J. Bonada, A. Loscos, O. Mayor, and H. Kenmochi, "Sample-based singing voice synthesizer using spectral models and source-filter decomposition," in *Third International Workshop on Models and Analysis of Vocal Emissions for Biomedical Applications*, 2003.

[3] H. Kenmochi and H. Ohshita, "Vocaloid-commercial singing synthesizer based on sample concatenation," in *Eighth Annual Conference of the International Speech Communication Association*, 2007.

[4] M. Nishimura, K. Hashimoto, K. Oura, Y. Nankaku, and K. Tokuda, "Singing voice synthesis based on deep neural networks." in *Interspeech*, 2016, pp. 2478–2482.

[5] J. Kim, H. Choi, J. Park, S. Kim, J. Kim, and M. Hahn, "Korean singing voice synthesis system based on an lstm recurrent neural network," in *INTERSPEECH 2018*. International Speech Communication Association, 2018.

[6] M. Blaauw and J. Bonada, "A neural parametric singing synthesizer modeling timbre and expression from natural songs," *Applied Sciences*, vol. 7, no. 12, p. 1313, 2017.

[7] H. Tachibana, K. Uenoyama, and S. Aihara, "Efficiently trainable text-to-speech system based on deep convolutional networks with guided attention," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 4784–4788.

[8] 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," *arXiv preprint arXiv:1703.10135*, 2017.

[9] A. Gibiansky, S. Arik, G. Diamos, J. Miller, K. Peng, W. Ping, J. Raiman, and Y. Zhou, "Deep voice 2: Multi-speaker neural text-to-speech," in *Advances in neural information processing systems*, 2017, pp. 2962–2970.

[10] R. Skerry-Ryan, E. Battenberg, Y. Xiao, Y. Wang, D. Stanton, J. Shor, R. J. Weiss, R. Clark, and R. A. Saurous, "Towards end-to-end prosody transfer for expressive speech synthesis with tacotron," *arXiv preprint arXiv:1803.09047*, 2018.

[11] Y. Wang, D. Stanton, Y. Zhang, R. Skerry-Ryan, E. Battenberg, J. Shor, Y. Xiao, F. Ren, Y. Jia, and R. A. Saurous, "Style tokens: Unsupervised style modeling, control and transfer in end-to-end speech synthesis," *arXiv preprint arXiv:1803.09017*, 2018.

[12] Y.-A. Chung, Y. Wang, W.-N. Hsu, Y. Zhang, and R. Skerry-Ryan, "Semi-supervised training for improving data efficiency in end-to-end speech synthesis," *arXiv preprint arXiv:1808.10128*, 2018.

[13] S. Li, S. Villette, P. Ramadas, and D. J. Sinder, "Speech bandwidth extension using generative adversarial networks," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 5029–5033.

[14] T. Miyato and M. Koyama, "cgans with projection discriminator," *arXiv preprint arXiv:1802.05637*, 2018.

[15] A. Brock, J. Donahue, and K. Simonyan, "Large scale gan training for high fidelity natural image synthesis," *arXiv preprint arXiv:1809.11096*, 2018.

[16] D. Griffin and J. Lim, "Signal estimation from modified short-time fourier transform," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 32, no. 2, pp. 236–243, 1984.

[17] A. v. d. Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, and K. Kavukcuoglu, "Wavenet: A generative model for raw audio," *arXiv preprint arXiv:1609.03499*, 2016.

[18] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: a simple way to prevent neural networks from overfitting," *The Journal of Machine Learning Research*, vol. 15, no. 1, pp. 1929–1958, 2014.

[19] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. C. Courville, "Improved training of wasserstein gans," in *Advances in Neural Information Processing Systems*, 2017, pp. 5767–5777.

[20] T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida, "Spectral normalization for generative adversarial networks," *arXiv preprint arXiv:1802.05957*, 2018.

[21] L. Mescheder, A. Geiger, and S. Nowozin, "Which training methods for gans do actually converge?" *arXiv preprint arXiv:1801.04406*, 2018.

[22] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in neural information processing systems*, 2014, pp. 2672–2680.

[23] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[24] 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*, 2010, pp. 249–256.

[25] 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.