# Speaking Speed Control of End-to-End Speech Synthesis using Sentence-Level Conditioning

*Jae-Sung Bae, Hanbin Bae, Young-Sun Joo, Junmo Lee, Gyeong-Hoon Lee, Hoon-Young Cho*

Speech AI Lab, NCSOFT, Republic of Korea

{jaesungbae, bhb0722, ysjoo555, ljun4121, ghlee0304, hycho}@ncsoft.com

## Abstract

This paper proposes a controllable end-to-end text-to-speech (TTS) system to control the speaking speed (speed-controllable TTS; SCTTS) of synthesized speech with sentence-level speaking-rate value as an additional input. The speaking-rate value, the ratio of the number of input phonemes to the length of input speech, is adopted in the proposed system to control the speaking speed. Furthermore, the proposed SCTTS system can control the speaking speed while retaining other speech attributes, such as the pitch, by adopting the global style token-based style encoder. The proposed SCTTS does not require any additional well-trained model or an external speech database to extract phoneme-level duration information and can be trained in an end-to-end manner. In addition, our listening tests on fast-, normal-, and slow-speed speech showed that the SCTTS can generate more natural speech than other phoneme duration control approaches which increase or decrease duration at the same rate for the entire sentence, especially in the case of slow-speed speech.

**Index Terms**: Text-to-speech, speed-controllable text-to-speech, speaking rate

## 1. Introduction

In the past few years, the naturalness of synthesized speech has been significantly improved owing to the advancement in neural text-to-speech (TTS) methods [1–4]. Speech varies in expressions; however, these models only focus on the generation of narrative-style speech. Therefore, many researches have been recently proposed to control the prosody and speaking speed of the synthesized speech in a TTS system [5–10]. This paper focuses on the control of speaking speed that is essential for real scenario because the speaking speed must vary depending on the context or situation.

In [5–8], various speech attributes, such as pitch, prosody, and speaking speed are extracted from a reference speech and the TTS system generates a similar style of speech to the reference speech. Because it is difficult to define every speech attribute objectively, the TTS system is trained in an unsupervised manner. However, because of this unsupervised learning mechanism, to successfully train and control speech attributes using these systems, the training database should contain a wide variety of styles of speech, and the amount of data in the database should be sufficient. Furthermore, it is very difficult to perfectly disentangle each speech attribute from speech (i.e., when we attempt to control the speaking speed, other speech attributes, such as pitch, may also be modified).

In [9, 10], neural TTS systems that control the phoneme-level speech duration have been proposed. Phoneme duration is additionally inputted to the TTS system [9], or the hidden states of the phoneme sequence are expanded, corresponding to the phoneme duration [10]. These systems, in the inference step, control the speaking speed by modifying the phoneme duration predicted by a duration predictor. However, these types of techniques have two constraints. First, they require an external well-trained automatic speech recognition (ASR) or TTS model to extract the ground-truth phoneme duration from the speech signal; however, to obtain these well-trained models, a database containing a large amount of extra speech data and domain knowledge are also required. Second, it is difficult to control the speaking speed naturally. These systems control the speaking speed by increasing or decreasing the predicted duration of each phoneme at the same rate for an entire sentence. However, when people are asked to speak quickly or slowly, their speaking speed is not consistent for the entire sentence; they control the speed by speaking specific phones quickly and others slowly [11]. In [12], the authors reported that the ratio of the duration of the vowels to that of the consonants increased as the speaking speed changed from normal to slow.

To overcome these constraints, we propose a speed-controllable TTS (SCTTS) system that adopts a sentence-level speaking rate (SR) as the input. The SR value is simply calculated as the ratio of the number of input phonemes to the length of input speech for each sentence. During training, by providing the SR value as the input, the proposed system learns to predict different alignments with the same text depending on desired speaking speeds. In the inference step, an average SR value obtained from a training database is used to generate a speech with normal speed, and the increased (decreased) SR value is used to generate fast (slow) speech. Furthermore, we improved the robustness of the speaking speed control of the proposed system by adopting a global style token (GST) [8]. In our pilot study, we found that the speed-control of a TTS system is affected by other speech attributes when the speech database contains various styles of speech samples (e.g., in addition to the speaking speed, other attributes such as pitch are changed). However, the GST-adopted SCTTS (SCTTS-GST) system can control the speaking speed while minimizing the changes in other speech attributes by disentangling it from multiple speech attributes included in the expressive speech. In the inference step, by providing a normal-style reference speech, together with the SR value, the SCTT-GST can control the speaking speed while retaining other speech attributes of a normal-style speech.

The key strengths of the proposed SCTTS system are as follows. First, the SCTTS system does not require extra labor-consuming well-trained models to extract the ground-truth phoneme duration and can be trained in an end-to-end manner. Second, the speed-controlled speech generated by the SCTTS system is more natural than that generated using other approaches of increasing or decreasing the phoneme duration at the same rate for the entire sentence. Third, the proposed SCTTS-GST system can control the speaking speed while retaining other speech attributes of a normal-style speech even for an expressive dataset.

## 2. Baseline Model

In this section, we present some background knowledge to understand our proposed method effectively.

### 2.1. End-to-End TTS Framework

As an end-to-end TTS framework, we used the deep convolutional TTS (DCTTS) system [4] with some modifications. Because the DCTTS system is fully convolutional, it has advantages of fast training speed and stable convergence. First, the text-to-mel spectrogram (T2M) network generates the coarse mel spectrogram, which is a down-sampled mel spectrogram in the time axis, from the input text. In the T2M network, the input phoneme (or character) embeddings and input acoustic features represented as the mel spectrogram are encoded by a text and an audio encoder, respectively. The attention module aligns the text and audio encoding, and the audio decoder predicts the coarse-mel spectrogram from the attention module output. Second, the post-processing network (PostNet) transforms the predicted coarse mel spectrogram into a mel spectrogram that is up-sampled in the time axis. As we synthesize speech using a neural vocoder for high-quality speech, the spectrogram super-resolution network (SSRN), which converts a coarse mel spectrogram into a spectrogram in [4], is replaced by the Post-Net. Finally, a neural vocoder synthesizes the speech waveform by inputting the predicted mel spectrogram.

Unlike the original DCTTS system, which trains T2M network and SSRN independently, we jointly train the T2M network and the PostNet in an end-to-end manner using the following loss function:

$$\mathcal{L} = \alpha \mathcal{L}_{spec}(\hat{c}|c) + \mathcal{L}_{spec}(\hat{m}|m), \quad (1)$$

where the $\hat{c}$ and $c$ denote the predicted and ground-truth coarse mel spectrograms, and $\hat{m}$ and $m$ are the predicted and ground-truth mel spectrograms. $\alpha$ is the weight of the coarse mel spectrogram reconstruction loss. The spectrogram reconstruction loss function, $\mathcal{L}_{spec}$, is defined as the summation of L1 loss and binary divergence $D_{bin}$ as in [4].

### 2.2. GST-based Style Encoder

The GST-based style encoder extracts the speaking style from the reference speech as a style embedding. First, the input reference speech signal is encoded by the reference encoder that was proposed in [5]. Then, the attention module generates the weights between the reference embedding and the GSTs, which represent their similarities. Finally, the weighted summation of GSTs forms the style embedding. In the training stage, the GSTs are randomly initialized and learned to contain the speech styles that appear dominant in the training dataset in an unsupervised manner. In the inference stage, the styles of synthesized speech can be controlled by the weights of GSTs. They can either be obtained by inputting the reference speech signal or be given directly.

## 3. Proposed Method

### 3.1. Speed-Controllable TTS Model

To control the speaking speed of the synthesized speech, we proposed the SCTTS system, which applies the sentence-level SR value to the end-to-end TTS system as an additional input. The architecture of the SCTTS system is detailed in Figure 1.

In the training stage, the text sequence, coarse mel spectrogram, and computed SR value are given as inputs. The SR value
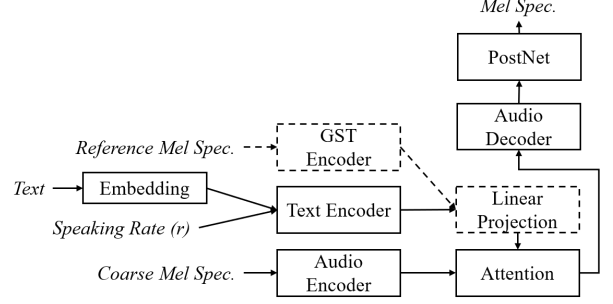


Figure 1: *Architecture of our proposed SCTTS system. Dotted-line components are added for the SCTTS-GST.*

is replicated to a text-embedding sequence length, concatenated with it, and then given as an input to the text encoder. In the inference stage, text sequences and the desired SR value are given, and the mel spectrogram with corresponding speaking speed is predicted.

### 3.2. Sentence-level SR Calculation

There are several methods of computing the SR value [13–15]. We used the inverse mean duration (IMD) [14] as the SR measurement. For each sentence, when the mel spectrogram with a length of $T$ and the text sequences with the number of $P$ are given, the SR value can be defined as follows:

$$r = \lambda \frac{P}{T}, \quad (2)$$

where $\lambda$ is a SR scaling factor. Please note that the silence regions are removed.

### 3.3. Disentanglement of Speaking Speed and Other Speech Attributes

We observe that multiple speech attributes are correlated when the speaking style of the speech dataset has significant variations. Figure 2 depicts the relationship between fundamental frequency ($F_0$) and SR of our neutral and highly expressive speech datasets; the speech dataset is detailed in section 4.1. In Figure 2, in our neutral dataset, the $F_0$ remains almost unchanged as the speaking speed changes, however, in our expressive dataset, it changes significantly. This type of correlation among speech attributes ($F_0$ and SR in Figure 2) affects speed control in the SCTTS system.

To disentangle and control the speaking speed from other speech attributes, we adopted a GST-based style encoder [8]. Because a GST-based style encoder learns style features in an unsupervised manner, in the SCTTS-GST system, it is expected to learn style features other than the speaking speed, which is provided as the SR value. The additional components of the SCTTS-GST system are indicated by dotted lines in Figure 1. The style embedding, which is the output of the GST-based style encoder, is concatenated with the text encoder output and projected through a linear layer. The other processes and loss function are the same as those of the basic SCTTS system.

## 4. Experiments and Results

We experimentally evaluated the following three aspects of the SCTTS system: 1) controlling the speaking speed, 2) disentangling the speaking speed from other speech attributes using
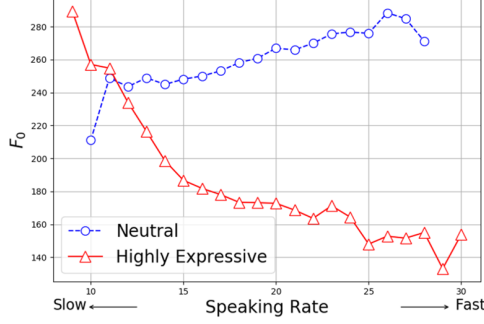
Figure 2: *$F_0$ vs. SR value of (red triangle) neutral-speaker and (blue circle) highly expressive-speaker training dataset.*
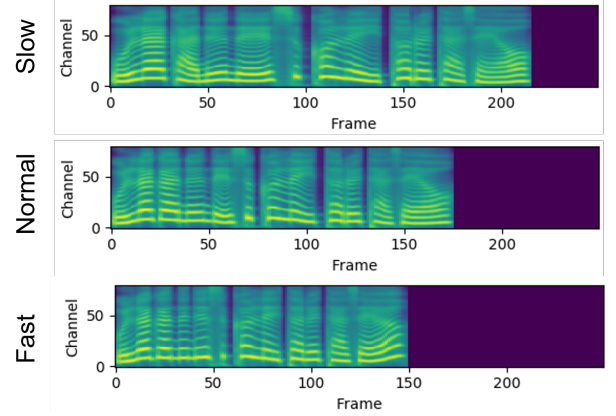


Figure 3: *Mel spectrograms of the synthesized speech with different speaking speeds. These were generated from the SCTTS system trained with the neutral speech dataset.*

GSTs, and 3) achieving naturalness of the synthesized speech with various speeds. Audio samples can be found online[1].

## 4.1. Datasets

For the experiments, we used two Korean speech datasets; one neutral and one highly expressive. The neutral speech dataset, KSS Dataset [16], was recorded by a professional female voice actress. The total amount of speech is about 8.7 h. The expressive speech dataset, an internal dataset, was recorded by a professional male voice actor. This dataset consists of speech samples with various expressions (neutral, excited, shouting, sigh, etc). The total amount of speech is 11 h. Both datasets were recorded in a studio environment and the frequency sampling rate is 22050 Hz. For testing, 1% of each dataset was randomly selected. We used phoneme as the text input and the mel spectrogram with 80 bins computed with an FFT size of 1024, a hop size of 256, and a window size of 1024, as the acoustic feature.

## 4.2. Model Setup

The architectures of the text encoder, audio encoder, and audio decoder were the same as in [4]. The architecture of Post-Net was the same as that of the SSRN in [4], but the output frequency bin was set to 80 because the PostNet predicts mel spectrogram. The SR scaling factor, $\lambda$, was set as 100. The loss weight, $\alpha$, was set as 1 before 50 k training steps; it decreased linearly and became 0 at the 200 k step. We did not use the guided attention loss in [4] because the system was trained stably and fast even without the guided attention loss. The Adam optimizer [17] was used to optimize the network with an initial learning rate of 0.001. For the GST-based style encoder, the same network architecture with ten style tokens and a multi-head attention module with four heads as in [8] were used. For a neural vocoder, WaveGlow [18] was used, and it was trained using a database containing approximately 50 h of speech recorded by four speakers.

## 4.3. Speaking Speed Control

Figure 3 shows the synthesized speech samples with various speaking speeds generated by the SCTTS system trained with the neutral speech dataset. We controlled the speed of the speech by adjusting the SR value. The average SR value representing the normal speed was obtained from the training set. As depicted in Figure 3, the SCTTS system naturally controlled the speaking speed while almost maintaining other speech at-

---

[1] https://nc-ai.github.io/speech/publications/speed-controllable-tts

tributes such as pitch, without leading to any distortions such as trembling or breaks in the mel spectrogram.

## 4.4. Disentanglement of Speaking Speed and Other Speech Attributes

In this sub-section, we evaluated the proposed system in the aspect of disentanglement between the speaking speed and other speech attributes. Figure 4 shows examples of speech samples with various speeds synthesized by the SCTTS and the SCTTS-GST systems. Both the systems were trained using the expressive speech dataset. For the SCTTS-GST system, we used two different reference speech, one with normal style (non-expressive; average $F_0$) and the other with excited style (high $F_0$). As illustrated in Figure 4, in the mel spectrograms of the synthesized speech from the SCTTS-GST system with normal and excited reference speech, $F_0$ remains unchanged as the speaking speed changes, whereas in the SCTTS system, $F_0$ changes depending on the speaking speed.

Figure 5 illustrates the average $F_0$ of the synthesized speech according to the SR values. For the test set, we synthesized speech samples for the same text by changing the SR value. As depicted in Figure 5, in the SCTTS, the $F_0$ was changed significantly according to the SR value. However, in the SCTTS-GST, the $F_0$ was nearly maintained regardless of the SR value (SCTTS-GST (N)) or changed much less (SCTTS-GST (H)) than the SCTTS. Because the speaking speed and high-$F_0$ are more strongly correlated in the training dataset (Figure 2).

## 4.5. Naturalness of the Speed-Controlled Speech

We evaluated the naturalness of the proposed system and compared it with phoneme-level speech duration control TTS systems. We compared the SCTTS and the SCTTS-GST systems with the phonemic-level duration controllable TTS (PDC-TTS) system [9] and the FastSpeech [10]. Please note that the same WaveGlow neural vocoder [18] was used in all the systems for high-quality speech generation.

### 4.5.1. Comparison Model Setup

**PDC-TTS:** For a fair comparison, we implemented the TTS part of the PDC-TTS by replacing it with DCTTS instead of Tacotron [2]. Accordingly, we concatenated the duration
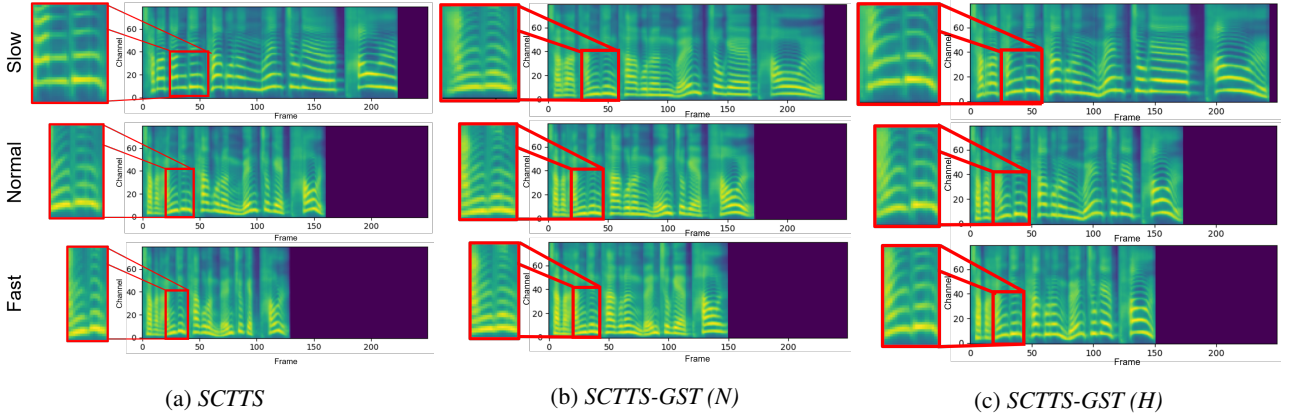
(a) *SCTTS*

(b) *SCTTS-GST (N)*

(c) *SCTTS-GST (H)*

Figure 4: *Mel spectrograms of the synthesized speech from the SCTTS and SCTTS-GST systems with different speaking speeds. "N" and "H" in the SCTTS-GST represent the normal- and high-$F_0$ reference speech cases, respectively. In (a), $F_0$ increased as the SR value decreased, depending on the speaker characteristics; however, in (b) and (c), by applying the GST technique, we successfully disentangled a person's speaking speed from $F_0$.*
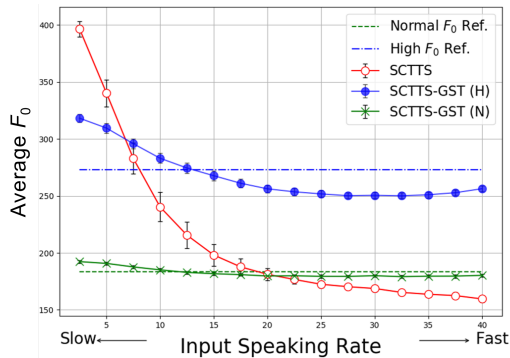


Figure 5: *Average $F_0$ vs. input SR value for the synthesized speech of the test set with 95% confidence interval (CI).*



Figure 6: *MOS test results on the naturalness with 95% CI.*

embedding, defined in [9], to the text encoder output, and fitted the output dimension through a linear projection. Because the phoneme duration model (i.e., duration predictor), which is essential in the speech-synthesis phase, is not mentioned in [9], we adopted the phoneme duration model in [19], except that the $F_0$ feature is removed from the input and output of the model. To train the duration model, the ground-truth phoneme duration was extracted from the speech database using a Kaldi-based ASR model [20]; the ASR model was pre-trained using an external ASR speech database. The trained phoneme duration model has a mean absolute error of 21.42 ms for the test set.

**FastSpeech:** We used open source implementation in the ESP-net framework [21] for training transformer TTS [22] and Fast-Speech. The pre-trained transformer-based TTS model was used to extract the ground-truth phoneme duration information.

*4.5.2. Subjective Evaluation*

We conducted MOS tests on the naturalness of the fast-, normal-, and slow-speed speech samples synthesized by each system. The length of the fast- and slow-speed speech were set as 70% and 150% of the length of the normal-speed speech, respectively. Each system was trained with the expressive dataset. For each system, fifteen speech samples were generated. A total of
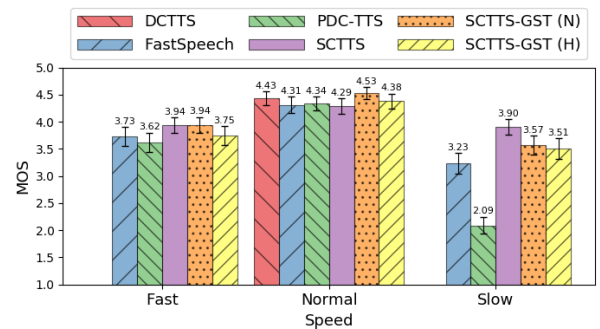
20 native Koreans participated. They were asked to score the naturalness of the synthesized-speech samples from 1 to 5. The naturalness of the speech synthesized by the DCTTS system is also compared for the normal-speed speech.

The MOS test results are presented in Figure 6. The performances of the proposed SCTTS and SCTTS-GST systems were comparable to or better than those of other approaches. Furthermore, at a slow speed, the proposed systems outperform other approaches drastically. Because the proposed systems generate slow-speed speech by lengthening the specific phoneme, word, or silence longer than the others, whereas other systems synthesize slow-speed speech by equally lengthening. The unnaturalness of speech generated from other approaches was more noticeable with the slow-speed speech than the fast-speed speech.

## 5. Conclusions

This paper proposed a speed-controllable TTS system that can generate a speech with various speaking speeds without the use of any pre-trained models and can be trained in an end-to-end manner. The proposed SCTTS system can be integrated with the GST and the speaking speed can be disentangled from other speech attributes such as pitch. The proposed SCTTS system outperformed the existing phoneme duration-control TTS systems in MOS listening tests, especially when the speed of the synthesized speech was slow.

# 6. References

[1] Y. Wang, R. J. Skerry-Ryan, D. Stanton, Y. Wu, R. J. Weiss, N. Jaitly, Z. Yang, Y. Xiao, Z. Chen, S. Bengio, Q. Le, Y. Agiomyrgiannakis, R. Clark, and R. A. Saurous, "Tacotron: Towards End-to-End Speech Synthesis," in *Proc. Interspeech*, 2017, pp. 4006–4010.

[2] J. Shen, R. Pang, R. J. Weiss, M. Schuster, N. Jaitly, Z. Yang, Z. Chen, Y. Zhang, Y. Wang, R. Skerrv-Ryan, R. A. Saurous, Y. Agiomvrgiannakis, and Y. Wu, "Natural TTS Synthesis by Conditioning Wavenet on MEL Spectrogram Predictions," in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2018, pp. 4779–4783.

[3] W. Ping, K. Peng, A. Gibiansky, S. Arik, A. Kannan, S. Narang, J. Raiman, and J. Miller, "Deep voice 3: 2000-speaker neural text-to-speech," in *Proc. International Conference on Learning Representations*, 2018.

[4] H. Tachibana, K. Uenoyama, and S. Aihara, "Efficiently Trainable Text-to-Speech System Based on Deep Convolutional Networks with Guided Attention," in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2018, pp. 4784–4788.

[5] R. Skerry-Ryan, E. Battenberg, Y. Xiao, Y. Wang, D. Stanton, J. Shor, R. Weiss, R. Clark, and R. Saurous, "Towards End-to-End Prosody Transfer for Expressive Speech Synthesis with Tacotron," in *Proc. of the 35th International Conference on Machine Learning*, 2018, pp. 4693–4702.

[6] V. Klimkov, S. Ronanki, J. Rohnke, and T. Drugman, "Fine-Grained Robust Prosody Transfer for Single-Speaker Neural Text-To-Speech," in *Proc. Interspeech*, 2019, pp. 4440–4444.

[7] Y. Lee and T. Kim, "Robust and Fine-grained Prosody Control of End-to-end Speech Synthesis," in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 5911–5915.

[8] 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," in *Proc. of the 35th International Conference on Machine Learning*, vol. 80, 2018, pp. 5180–5189.

[9] J. Park, K. Han, Y. Jeong, and S. W. Lee, "Phonemic-level Duration Control Using Attention Alignment for Natural Speech Synthesis," in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 5896–5900.

[10] Y. Ren, Y. Ruan, X. Tan, T. Qin, S. Zhao, Z. Zhao, and T. Liu, "Fastspeech: Fast, robust and controllable text to speech," in *Proc. Advances in Neural Information Processing Systems*, 2019, pp. 3165–3174.

[11] F. Martinez, D. Tapias, J. Alvarez, and P. Leon, "Characteristics of slow, average and fast speech and their effects in large vocabulary continuous speech recognition," in *Fifth European Conference on Speech Communication and Technology*, 1997.

[12] H. Kuwabara, "Acoustic properties of phonemes in continuous speech for different speaking rate," in *Proc. of Fourth International Conference on Spoken Language Processing. (ICSLP)*, vol. 4. IEEE, 1996, pp. 2435–2438.

[13] M. A. Siegler and R. M. Stern, "On the effects of speech rate in large vocabulary speech recognition systems," in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, vol. 1, 1995, pp. 612–615.

[14] N. Mirghafori, E. Foster, and N. Morgan, "Fast speakers in large vocabulary continuous speech recognition: analysis & antidotes," in *Fourth European Conference on Speech Communication and Technology*, 1995, p. 491494.

[15] F. Martinez, D. Tapias, J. Alvarez, and P. Leon, "Characteristics of slow, average and fast speech and their effects in large vocabulary continuous speech recognition," in *Fifth European Conference on Speech Communication and Technology*, 1997.

[16] K. Park, "KSS Dataset: Korean Single speaker Speech Dataset," https://kaggle.com/bryanpark/korean-single-speaker-speech-dataset, 2018.

[17] D. P. Kingma and J. Ba, "Adam: A Method for Stochastic Optimization," *CoRR*, vol. abs/1412.6980, 2014. [Online]. Available: http://arxiv.org/abs/1412.6980

[18] R. Prenger, R. Valle, and B. Catanzaro, "WaveGlow: A Flow-based Generative Network for Speech Synthesis," in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 3617–3621.

[19] S. Ö. Arik, M. Chrzanowski, A. Coates, G. Diamos, A. Gibiansky, Y. Kang, X. Li, J. Miller, A. Ng, J. Raiman *et al.*, "Deep voice: Real-time neural text-to-speech," in *Proc. International Conference on Machine Learning (ICML)*, 2017, pp. 195–204.

[20] D. Povey, A. Ghoshal, G. Boulianne, L. Burget, O. Glembek, N. Goel, M. Hannemann, P. Motlicek, Y. Qian, P. Schwarz *et al.*, "The Kaldi speech recognition toolkit," in *IEEE 2011 workshop on automatic speech recognition and understanding*, no. CONF. IEEE Signal Processing Society, 2011.

[21] S. Watanabe, T. Hori, S. Karita, T. Hayashi, J. Nishitoba, Y. Unno, N. Enrique Yalta Soplin, J. Heymann, M. Wiesner, N. Chen, A. Renduchintala, and T. Ochiai, "ESPnet: End-to-End Speech Processing Toolkit," in *Proc. Interspeech*, 2018, pp. 2207–2211.

[22] N. Li, S. Liu, Y. Liu, S. Zhao, and M. Liu, "Neural speech synthesis with transformer network," in *Proc. of the AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 6706–6713.