# Learning to Speak Fluently in a Foreign Language:
# Multilingual Speech Synthesis and Cross-Language Voice Cloning

*Yu Zhang, Ron J. Weiss, Heiga Zen, Yonghui Wu, Zhifeng Chen, RJ Skerry-Ryan, Ye Jia,*
*Andrew Rosenberg, Bhuvana Ramabhadran*

Google

{ngyuzh, ronw}@google.com

## Abstract

We present a multispeaker, multilingual text-to-speech (TTS) synthesis model based on Tacotron that is able to produce high quality speech in multiple languages. Moreover, the model is able to transfer voices across languages, e.g. synthesize fluent Spanish speech using an English speaker's voice, without training on any bilingual or parallel examples. Such transfer works across distantly related languages, e.g. English and Mandarin.

Critical to achieving this result are: 1. using a phonemic input representation to encourage sharing of model capacity across languages, and 2. incorporating an adversarial loss term to encourage the model to disentangle its representation of speaker identity (which is perfectly correlated with language in the training data) from the speech content. Further scaling up the model by training on multiple speakers of each language, and incorporating an autoencoding input to help stabilize attention during training, results in a model which can be used to consistently synthesize intelligible speech for training speakers in all languages seen during training, and in native or foreign accents.

**Index Terms**: speech synthesis, end-to-end, adversarial loss

## 1. Introduction

Recent end-to-end neural TTS models [1–3] have been extended to enable control of speaker identity [4–7] as well as unlabelled speech attributes, e.g. prosody, by conditioning synthesis on latent representations [8–12] in addition to text. Extending such models to support multiple, unrelated languages is nontrivial when using language-dependent input representations or model components, especially when the amount of training data per language is imbalanced. For example, there is no overlap in the text representation between languages like Mandarin and English. Furthermore, recordings from bilingual speakers are expensive to collect. It is therefore most common for each speaker in the training set to speak only one language, so speaker identity is perfectly correlated with language. This makes it difficult to transfer voices across different languages, a desirable feature when the number of available training voices for a particular language is small. Moreover, for languages with borrowed or shared words, such as proper nouns in Spanish (ES) and English (EN), pronunciations of the same text might be different. This adds more ambiguity when a naively trained model sometimes generates accented speech for a particular speaker.

Zen et al. proposed a speaker and language factorization for HMM-based parametric TTS system [13], aiming to transfer a voice from one language to others. [14] proposed a multilingual parametric neural TTS system, which used a unified input representation and shared parameters across languages, however the voices used for each language were disjoint. [15] described a similar bilingual Chinese and English neural TTS system trained on speech from a bilingual speaker, allowing it to synthesize speech
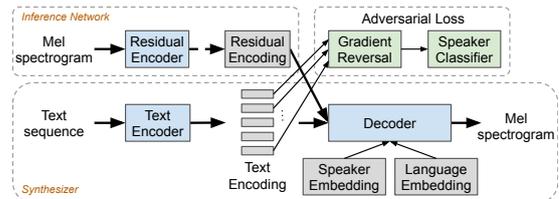


Figure 1: *Overview of the components of the proposed model. Dashed lines denote sampling via reparameterization [21] during training. The prior mean is always use during inference.*

in both languages using the same voice. [16] studied learning pronunciation from a bilingual TTS model. Most recently, [17] presented a multilingual neural TTS model which supports voice cloning across English, Spanish, and German. It used language-specific text and speaker encoders, and incorporated a secondary fine-tuning step to optimize a speaker identity-preserving loss, ensuring that the model could output a consistent voice regardless of language. We also note that the sound quality is not on par with recent neural TTS systems, potentially because of its use of the WORLD vocoder [18] for waveform synthesis.

Our work is most similar to [19], which describes a multilingual TTS model based on Tacotron 2 [20] which uses a Unicode encoding "byte" input representation to train a model on one speaker of each of English, Spanish, and Mandarin. In this paper, we evaluate different input representations, scale up the number of training speakers for each language, and extend the model to support cross-lingual voice cloning. The model is trained in a single stage, with no language-specific components, and obtains naturalness on par with baseline monolingual models. Our contributions include: (1) Evaluating the effect of using different text input representations in a multilingual TTS model. (2) Introducing a per-input token speaker-adversarial loss to enable cross-lingual voice transfer when only one training speaker is available for each language. (3) Incorporating an explicit language embedding to the input, which enables moderate control of speech accent, independent of speaker identity, when the training data contains multiple speakers per language.

We evaluate the contribution of each component, and demonstrate the proposed model's ability to disentangle speakers from languages and consistently synthesize high quality speech for all speakers, despite the perfect correlation to the original language in the training data.

## 2. Model Structure

We base our multilingual TTS model on Tacotron 2 [20], which uses an attention-based sequence-to-sequence model to generate a sequence of log-mel spectrogram frames based on an input text sequence. The architecture is illustrated in Figure 1. It

augments the base Tacotron 2 model with additional speaker and, optionally, language embedding inputs (bottom right), an adversarially-trained speaker classifier (top right), and a variational autoencoder-style residual encoder (top left) which conditions the decoder on a latent embedding computed from the target spectrogram during training (top left). Finally, similar to Tacotron 2, we separately train a WaveRNN [22] neural vocoder.

## 2.1. Input representations

End-to-end TTS models have typically used character [2] or phoneme [8, 23] input representations, or hybrids between them [24, 25]. Recently, [19] proposed using inputs derived from the UTF-8 byte encoding in multilingual settings. We evaluate the effects of using these representations for multilingual TTS.

### 2.1.1. Characters / Graphemes

Embeddings corresponding to each character or grapheme are the default inputs for end-to-end TTS models [2, 20, 23], requiring the model to implicitly learn how to pronounce input words (i.e. grapheme-to-phoneme conversion [26]) as part of the synthesis task. Extending a grapheme-based input vocabulary to a multilingual setting is straightforward, by simply concatenating grapheme sets in the training corpus for each language. This can grow quickly for languages with large alphabets, e.g. our Mandarin vocabulary contains over 4.5k tokens. We simply concatenate all graphemes appearing in the training corpus, leading to a total of 4,619 tokens. Equivalent graphemes are shared across languages. During inference all previously unseen characters are mapped to a special out-of-vocabulary (OOV) symbol.

### 2.1.2. UTF-8 Encoded Bytes

Following [19] we experiment with an input representation based on the UTF-8 text encoding, which uses 256 possible values as each input token where the mapping from graphemes to bytes is language-dependent. For languages with single-byte characters (e.g., English), this representation is equivalent to the grapheme representation. However, for languages with multi-byte characters (such as Mandarin) the TTS model must learn to attend to a consistent sequence of bytes to correctly generate the corresponding speech. On the other hand, using a UTF-8 byte representation may promote sharing of representations between languages due to the smaller number of input tokens.

### 2.1.3. Phonemes

Using phoneme inputs simplifies the TTS task, as the model no longer needs to learn complicated pronunciation rules for languages such as English. Similar to our grapheme-based model, equivalent phonemes are shared across languages. We concatenate all possible phoneme symbols, for a total of 88 tokens.

To support Mandarin, we include tone information by learning phoneme-independent embeddings for each of the 4 possible tones, and broadcast each tone embedding to all phoneme embeddings inside the corresponding syllable. For English and Spanish, tone embeddings are replaced by stress embeddings which include primary and secondary stresses. A special symbol is used when there is no tone or stress.

## 2.2. Residual encoder

Following [12], we augment the TTS model by incorporating a variational autoencoder-like *residual encoder* which encodes the latent factors in the training audio, e.g. prosody or background noise, which is not well-explained by the conditioning inputs: the text representation, speaker, and language embeddings. We follow the structure from [12], except we use a standard single Gaussian prior distribution and reduce the latent dimension to 16. In our experiments, we observe that feeding in the prior mean (all zeros) during inference, significantly improves stability of cross-lingual speaker transfer and leads to improved naturalness as shown by MOS evaluations in Section 3.4.

## 2.3. Adversarial training

One of the challenges for multilingual TTS is data sparsity, where some languages may only have training data for a few speakers. In the extreme case where there is only one speaker per language in the training data, the speaker identity is essentially the same as the language ID. To encourage the model to learn disentangled representations of the text and speaker identity, we proactively discourage the text encoding $t_s$ from also capturing speaker information. We employ domain adversarial training [27] to encourage $t_i$ to encode text in a speaker-independent manner by introducing a speaker classifier based on the text encoding and a gradient reversal layer. Note that the speaker classifier is optimized with a different objective than the rest of the model: $\mathcal{L}_{\text{speaker}}(\psi_s; t_i) = \sum_i^N \log p(s_i \mid t_i)$, where $s_i$ is the speaker label and $\psi_s$ are the parameters for speaker classifier. To train the full model, we insert a gradient reversal layer [27] prior to this speaker classifier, which scales the gradient by $-\lambda$. Following [28], we also explore inserting another adversarial layer on top of the variational autoencoder to encourage it to learn speaker-independent representations. However, we found that this layer has no effect after decreasing the latent space dimension.

We impose this adversarial loss separately on each element of the encoded text sequence, in order to encourage the model to learn a speaker- and language-independent text embedding space. In contrast to [28], which disentangled speaker identity from background noise, some input tokens are highly language-dependent which can lead to unstable adversarial classifier gradients. We address this by clipping gradients computed at the reversal layer to limit the impact of such outliers.

# 3. Experiments

We train models using a proprietary dataset composed of high quality speech in three languages: (1) 385 hours of English (EN) from 84 professional voice actors with accents from the United States, Great Britain, Australia, and Singapore; (2) 97 hours of Spanish (ES) from 3 female speakers include Castilian and US Spanish; (3) 68 hours of Mandarin (CN) from 5 speakers.

## 3.1. Model and training setup

The synthesizer network uses the Tacotron 2 architecture [20], with additional inputs consisting of learned speaker (64-dim) and language embeddings (3-dim), concatenated and passed to the decoder at each step. The generated speech is represented as a sequence of 128-dim log-mel spectrogram frames, computed from 50ms windows shifted by 12.5ms.

The variational residual encoder architecture closely follows the attribute encoder in [12]. It maps a variable length mel spectrogram to two vectors parameterizing the mean and log variance of the Gaussian posterior. The speaker classifiers are fully-connected networks with one 256 unit hidden layer followed by a softmax predicting the speaker identity. The synthesizer and speaker classifier are trained with weight 1.0 and 0.02 respectively. As described in the previous section we apply

Table 1: *Speaker similarity Mean Opinion Score (MOS) comparing ground truth audio from speakers of different languages. Raters are native speakers of the target language.*

| Source Language | Target Language | | |
|---|---|---|---|
| | EN | ES | CN |
| EN | 4.40±0.07 | 1.72±0.15 | 1.80±0.08 |
| ES | 1.49±0.06 | 4.39±0.06 | 2.14±0.09 |
| CN | 1.32±0.06 | 2.06±0.09 | 3.51±0.12 |

gradient clipping with factor 0.5 to the gradient reversal layer.

The entire model is trained jointly with a batch size of 256, using the Adam optimizer configured with an initial learning rate of $10^{-3}$, and an exponential decay that halves the learning rate every 12.5k steps, starting at 50k steps.

Waveforms are synthesized using a WaveRNN [22] vocoder which generates 16-bit signals sampled at 24 kHz conditioned on spectrograms predicted by the TTS model. We synthesize 100 samples per model, and have each one rated by 6 raters.

### 3.2. Evaluation

To evaluate synthesized speech, we rely on crowdsourced Mean Opinion Score (MOS) evaluations of *speech naturalness* via subjective listening tests. Ratings follow the Absolute Category Rating scale, with scores from 1 to 5 in 0.5 point increments.

For cross-language voice cloning, we also evaluate whether the synthesized speech resembles the identity of the reference speaker by pairing each synthesized utterance with a reference utterance from the same speaker for subjective MOS evaluation of *speaker similarity*, as in [5]. Although rater instructions explicitly asked for the content to be ignored, note that this similarity evaluation is more challenging than the one in [5] because the reference and target examples are spoken in different languages, and raters are not bilingual. We found that low fidelity audio tended to result in high variance similarity MOS so we always use WaveRNN outputs.[1]

For each language, we chose one speaker to use for similarity tests. As shown in Table 1, the EN speaker is found to be dissimilar to the ES and CN speakers (MOS below 2.0), while the ES and CN speakers are slightly similar (MOS around 2.0). The CN speaker has more natural variability compared to EN and ES, leading to a lower self similarity. The scores are consistent when EN and CN raters evaluate the same EN and CN test set. The observation is consistent with [29]: raters are able to discriminate between speakers across languages. However, when rating synthetic speech, we observed that English speaking raters often considered "heavy accented" synthetic CN speech to sound more similar to the target EN speaker, compared to more fluent speech from the same speaker. This indicates that accent and speaker identity are not fully disentangled. We encourage readers to listen to samples on the companion webpage.[2]

### 3.3. Comparing input representations

We first build and evaluate models comparing the performance of different text input representations. For all three languages, byte-based models always use a 256-dim softmax output. Monolingual character and phoneme models each use a different input

Table 2: *Naturalness MOS of monolingual and multilingual models synthesizing speech of in different languages.*

| Model | Input | Language | | |
|---|---|---|---|---|
| | | EN | ES | CN |
| Ground truth | | 4.60±0.05 | 4.37±0.06 | 4.42±0.06 |
| Monolingual | char | 4.24±0.12 | 4.21±0.11 | 3.48±0.11 |
| | phone | 4.59±0.06 | 4.39±0.04 | 4.16±0.08 |
| Multilingual 1EN 1ES 1CN | byte | 4.23±0.14 | 4.23±0.10 | 3.42±0.12 |
| | char | 3.94±0.15 | 4.33±0.09 | 3.63±0.10 |
| | phone | 4.34±0.09 | 4.41±0.05 | 4.06±0.10 |
| Multilingual 84EN 3ES 5CN | byte | 4.11±0.14 | 4.21±0.12 | 3.67±0.12 |
| | char | 4.26±0.13 | 4.23±0.11 | 3.46±0.11 |
| | phone | 4.37±0.12 | 4.37±0.04 | 4.09±0.10 |

Table 3: *Naturalness and speaker similarity MOS of cross-language voice cloning of an EN source speaker. Models which use different input representations are compared, with and without the speaker-adversarial loss. fail: raters complained that too many utterances were spoken in the wrong language.*

| Input | ES target | | CN target | |
|---|---|---|---|---|
| | Naturalness | Similarity | Naturalness | Similarity |
| char | 2.62±0.10 | 4.25±0.09 | N/A | N/A |
| byte | 2.62±0.15 | 3.96±0.10 | N/A | N/A |
| with adversarial loss | | | | |
| byte | 2.34±0.10 | 4.23±0.09 | fail | 3.85±0.11 |
| phone | 3.20±0.09 | 4.15±0.10 | 2.75±0.12 | 3.60±0.09 |

vocabulary corresponding to the training language.

Table 2 compares monolingual and multilingual model performance using different input representations. For Mandarin, the phoneme-based model performs significantly better than char- or byte-based variants due to rare and OOV words. Compared to the monolingual system, multilingual phoneme-based systems have similar performance on ES and CN but are slightly worse on EN. CN has a larger gap to ground truth (top) due to unseen word segmentation (for simplicity, we didn't add word boundary during training). The multispeaker model (bottom) performs about the same as the single speaker per-language variant (middle). Overall, when using phoneme inputs all the languages obtain MOS scores above 4.0.

### 3.4. Cross-language voice cloning

We evaluate how well the multispeaker models can be used to clone a speaker's voice into a new language by simply passing in speaker embeddings corresponding to a different language from the input text. Table 3 shows voice cloning performance from an EN speaker in the most data-poor scenario (129 hours), where only a single speaker is available for each training language (1EN 1ES 1CN) without using the speaker-adversarial loss. Using byte inputs [3] it was possible to clone the EN speaker to ES with high similarity MOS, albeit with significantly reduced naturalness. However, cloning the EN voice to CN failed[4], as did cloning to ES and CN using phoneme inputs.

---

[1]Some raters gave low fidelity audio lower scores, treating "blurriness" as a property of the speaker. Others gave higher scores because they recognized such audio as synthetic and had lower expectations.

[3]Using character or byte inputs led to similar results.

[4]We didn't run listening tests because it was clear that synthesizing EN text using the CN speaker embedding didn't affect the model output.

Table 4: *Naturalness and speaker similarity MOS of cross-language voice cloning of the full multilingual model using phoneme inputs.*

| Source Language | Model | EN target | | ES target | | CN target | |
|---|---|---|---|---|---|---|---|
| | | Naturalness | Similarity | Naturalness | Similarity | Naturalness | Similarity |
| - | Ground truth (self-similarity) | 4.60±0.05 | 4.40±0.07 | 4.37±0.06 | 4.39±0.06 | 4.42±0.06 | 3.51±0.12 |
| EN | 84EN 3ES 5CN | 4.37±0.12 | 4.63±0.06 | 4.20±0.07 | 3.50±0.12 | 3.94±0.09 | 3.03±0.10 |
| | language ID fixed to EN | - | - | 3.68±0.07 | 4.06±0.09 | 3.09±0.09 | 3.20±0.09 |
| ES | 84EN 3ES 5CN | 4.28±0.10 | 3.24±0.09 | 4.37±0.04 | 4.01±0.07 | 3.85±0.09 | 2.93±0.12 |
| CN | 84EN 3ES 5CN | 4.49±0.08 | 2.46±0.10 | 4.56±0.08 | 2.48±0.09 | 4.09±0.10 | 3.45±0.12 |

Adding the adversarial speaker classifier enabled cross-language cloning of the EN speaker to CN with very high similarity MOS for both byte and phoneme models. However, naturalness MOS remains much lower than using the native speaker identity, with the naturalness listening test failing entirely in the CN case with byte inputs as a result of rater comments that the speech sounded like a foreign language. According to rater comments on the phoneme system, most of the degradation came from mismatched accent and pronunciation, not fidelity. CN raters commented that it sounded like "a foreigner speaking Chinese". More interestingly, few ES raters commented that "The voice does not sound robotic but instead sounds like an English native speaker who is learning to pronounce the words in Spanish." Based on these results, we only use phoneme inputs in the following experiments since this guarantees that pronunciations are correct and results in more fluent speech.

Table 4 evaluates voice cloning performance of the full multilingual model (84EN 3ES 5CN), which is trained on the full dataset with increased speaker coverage, and uses the speaker-adversarial loss and speaker/language embeddings. Incorporating the adversarial loss forces the text representation to be less language-specific, instead relying on the language embedding to capture language-dependent information. Across all language pairs, the model synthesizes speech in all voices with naturalness MOS above 3.85, demonstrating that increasing training speaker diversity improves generalization. In most cases synthesizing EN and ES speech (except EN-to-ES) approaches the ground truth scores. In contrast, naturalness of CN speech is consistently lower than the ground truth.

The high naturalness and similarity MOS scores in the top row of Table 4 indicate that the model is able to successfully transfer the EN voice to both ES and CN almost without accent. When consistently conditioning on the EN language embedding regardless of the target language (second row), the model produces more English accented ES and CN speech, which leads to lower naturalness but higher similarity MOS scores. Also see Figure 2 and the demo for accent transfer audio examples.

We see that cloning the CN voice to other languages (bottom row) has the lowest similarity MOS, although the scores are still much higher than different-speaker similarity MOS in the off-diagonals of Table 1 indicating that there is some degree of transfer. This is a consequence of the low speaker coverage of CN compared to EN in the training data, as well as the large distance between CN and other languages.

Finally, Table 5 demonstrates the importance of training using a variational residual encoder to stabilize the model output. Naturalness MOS decreases by 0.4 points for EN-to-CN cloning without the residual encoder (bottom row). In informal comparisons of the outputs of the two models we find that the model without the residual encoder tends to skip rare words or inserts

Table 5: *Effect of EN speaker cloning with no residual encoder.*

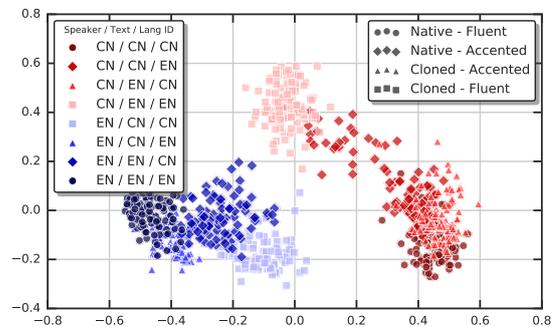| Model | Target Language | | |
|---|---|---|---|
| | EN | ES | CN |
| 84EN 3ES 5CN | 4.37±0.12 | 4.20±0.07 | 3.94±0.09 |
| - residual encoder | 4.38±0.10 | 4.11±0.06 | 3.52±0.11 |



Figure 2: *Visualizing the effect of voice cloning and accent control, using 2D PCA of speaker embeddings [30] computed from speech synthesized with different speaker, text, and language ID combinations. Embeddings cluster together (bottom left and right), implying high similarity, when the speaker's original language matches the language embedding, regardless of the text language. However, using language ID from the text (squares), modifying the speaker's accent to speak fluently, hurts similarity compared to the native language and accent (circles).*

unnatural pauses in the output speech. This indicates the VAE prior learns a mode which helps stabilize attention.

## 4. Conclusions

We describe extensions to the Tacotron 2 neural TTS model which allow training of a multilingual model trained only on monolingual speakers, which is able to synthesize high quality speech in three languages, and transfer training voices across languages. Furthermore, the model learns to speak foreign languages with moderate control of accent, and, as demonstrated on the companion webpage, has rudimentary support for code switching. In future work we plan to investigate methods for scaling up to leverage large amounts of low quality training data, and support many more speakers and languages.

## 5. Acknowledgements

# 6. References

[1] A. van den 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," *CoRR abs/1609.03499*, 2016.

[2] 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: A fully end-to-end text-to-speech synthesis model," *arXiv preprint*, 2017.

[3] S. Arik, G. Diamos, A. Gibiansky, 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 (NIPS)*, 2017.

[4] S. O. Arik, J. Chen, K. Peng, W. Ping, and Y. Zhou, "Neural voice cloning with a few samples," in *Advances in Neural Information Processing Systems*, 2018.

[5] Y. Jia, Y. Zhang, R. J. Weiss, Q. Wang, J. Shen, F. Ren, Z. Chen, P. Nguyen, R. Pang, I. L. Moreno, and Y. Wu, "Transfer learning from speaker verification to multispeaker text-to-speech synthesis," in *Advances in Neural Information Processing Systems*, 2018.

[6] E. Nachmani, A. Polyak, Y. Taigman, and L. Wolf, "Fitting new speakers based on a short untranscribed sample," in *International Conference on Machine Learning (ICML)*, 2018.

[7] Y. Chen, Y. Assael, B. Shillingford, D. Budden, S. Reed, H. Zen, Q. Wang, L. C. Cobo, A. Trask, B. Laurie *et al.*, "Sample efficient adaptive text-to-speech," *arXiv preprint arXiv:1809.10460*, 2018.

[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 *International Conference on Machine Learning (ICML)*, 2018.

[9] 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," in *International Conference on Machine Learning (ICML)*, 2018.

[10] K. Akuzawa, Y. Iwasawa, and Y. Matsuo, "Expressive speech synthesis via modeling expressions with variational autoencoder," in *Interspeech*, 2018.

[11] G. E. Henter, J. Lorenzo-Trueba, X. Wang, and J. Yamagishi, "Deep encoder-decoder models for unsupervised learning of controllable speech synthesis," *arXiv preprint arXiv:1807.11470*, 2018.

[12] W.-N. Hsu, Y. Zhang, R. J. Weiss, H. Zen, Y. Wu, Y. Wang, Y. Cao, Y. Jia, Z. Chen, J. Shen, P. Nguyen, and R. Pang, "Hierarchical generative modeling for controllable speech synthesis," in *ICLR*, 2019.

[13] H. Zen, N. Braunschweiler, S. Buchholz, M. Gales, K. Knill, S. Krstulović, and J. Latorre, "Statistical parametric speech synthesis based on speaker and language factorization," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 20, no. 6, pp. 1713–1724, 2012.

[14] B. Li and H. Zen, "Multi-language multi-speaker acoustic modeling for LSTM-RNN based statistical parametric speech synthesis," in *Proc. Interspeech*, 2016, pp. 2468–2472.

[15] H. Ming, Y. Lu, Z. Zhang, and M. Dong, "A light-weight method of building an LSTM-RNN-based bilingual TTS system," in *International Conference on Asian Language Processing*, 2017, pp. 201–205.

[16] Y. Lee and T. Kim, "Learning pronunciation from a foreign language in speech synthesis networks," *arXiv preprint arXiv:1811.09364*, 2018.

[17] E. Nachmani and L. Wolf, "Unsupervised polyglot text to speech," in *ICASSP*, 2019.

[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] B. Li, Y. Zhang, T. Sainath, Y. Wu, and W. Chan, "Bytes are all you need: End-to-end multilingual speech recognition and synthesis with bytes," in *ICASSP*, 2018.

[20] J. Shen, R. Pang, R. J. Weiss, M. Schuster, N. Jaitly, Z. Yang, Z. Chen, Y. Zhang, Y. Wang, R. Skerry-Ryan *et al.*, "Natural TTS synthesis by conditioning WaveNet on mel spectrogram predictions," in *ICASSP*, 2018.

[21] D. P. Kingma and M. Welling, "Auto-encoding variational Bayes," in *International Conference on Learning Representations (ICLR)*, 2014.

[22] N. Kalchbrenner, E. Elsen, K. Simonyan, S. Noury, N. Casagrande, E. Lockhart, F. Stimberg, A. van den Oord, S. Dieleman, and K. Kavukcuoglu, "Efficient neural audio synthesis," in *ICML*, 2018.

[23] J. Sotelo, S. Mehri, K. Kumar, J. F. Santos, K. Kastner, A. Courville, and Y. Bengio, "Char2wav: End-to-end speech synthesis," in *ICLR: Workshop*, 2017.

[24] W. Ping, K. Peng, A. Gibiansky, S. O. Arik, A. Kannan, S. Narang, J. Raiman, and J. Miller, "Deep Voice 3: Scaling text-to-speech with convolutional sequence learning," in *International Conference on Learning Representations (ICLR)*, 2018.

[25] K. Kastner, J. F. Santos, Y. Bengio, and A. C. Courville, "Representation mixing for TTS synthesis," *arXiv:1811.07240*, 2018.

[26] A. Van Den Bosch and W. Daelemans, "Data-oriented methods for grapheme-to-phoneme conversion," in *Proc. Association for Computational Linguistics*, 1993, pp. 45–53.

[27] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky, "Domain-adversarial training of neural networks," *The Journal of Machine Learning Research*, vol. 17, no. 1, pp. 2096–2030, 2016.

[28] W.-N. Hsu, Y. Zhang, R. J. Weiss, Y. an Chung, Y. Wang, Y. Wu, and J. Glass, "Disentangling correlated speaker and noise for speech synthesis via data augmentation and adversarial factorization," in *ICASSP*, 2019.

[29] M. Wester and H. Liang, "Cross-lingual speaker discrimination using natural and synthetic speech," in *Twelfth Annual Conference of the International Speech Communication Association*, 2011.

[30] L. Wan, Q. Wang, A. Papir, and I. L. Moreno, "Generalized end-to-end loss for speaker verification," in *Proc. ICASSP*, 2018.