



Counterfactual Activation Editing for Post-hoc Prosody and Mispronunciation Correction in TTS Models

Kywoon Lee¹, Artyom Stitsyuk¹, Gunu Jho², Inchul Hwang², Jaesik Choi^{1,3}

¹KAIST, South Korea

²Samsung Electronics, South Korea

³INEEJI, South Korea

{leekwoon, stitsyuk, jaesik.choi}@kaist.ac.kr, {gunu.jho, inc.hwang}@samsung.com

Abstract

Recent advances in Text-to-Speech (TTS) have significantly improved speech naturalness, increasing the demand for precise prosody control and mispronunciation correction. Existing approaches for prosody manipulation often depend on specialized modules or additional training, limiting their capacity for post-hoc adjustments. Similarly, traditional mispronunciation correction relies on grapheme-to-phoneme dictionaries, making it less practical in low-resource settings. We introduce Counterfactual Activation Editing, a model-agnostic method that manipulates internal representations in a pre-trained TTS model to achieve post-hoc control of prosody and pronunciation. Experimental results show that our method effectively adjusts prosodic features and corrects mispronunciations while preserving synthesis quality. This opens the door to inference-time refinement of TTS outputs without retraining, bridging the gap between pre-trained TTS models and editable speech synthesis.

Index Terms: speech synthesis, prosody control, pronunciation control

1. Introduction

Advancements in Text-to-Speech (TTS) models [1, 2, 3, 4, 5] and neural vocoders [6, 7, 8, 9] have made synthetic voices nearly indistinguishable from human speech. Consequently, more attention has been attracted by refining the expressiveness of synthetic voices, particularly through the accurate manipulation of prosodic features and pronunciation.

Traditionally, controlling prosody and correcting mispronunciations in TTS models have involved distinct approaches. For prosody control, one popular strategy has been learning a latent space, from which embeddings are sampled during inference to supplement the text with necessary prosodic information [10, 11, 12, 13]. However, due to its unsupervised nature, selecting the right embedding can be challenging. An alternative strategy involves adding specific layers to the TTS model that are trained to extract and manipulate this prosodic information directly from the input sequence [4, 5, 14, 15], but this requires extensive preprocessing. Both of these methods, however, lack straightforward applicability to already trained models, such as Tacotron 2 [2] or Transformer TTS [3], without undergoing a retraining process.

Similarly, the field of pronunciation correction has heavily relied on pronouncing dictionaries to convert grapheme to phoneme. While effective, this method struggles with scalability, precluding TTS models from being useful in low-resource languages. Diversifying pronunciations during training helps, but the skewed word distribution in natural language (Zipfian

distribution) [16] limits its utility. In general, TTS training datasets cannot match the comprehensive coverage of pronunciation dictionaries, highlighting the need for a more flexible and encompassing strategy.

In this paper, we introduce an *orthogonal* framework leveraging *Counterfactual Activation Editing* (CAE) to directly manipulate the prosody and pronunciation in pretrained TTS models. Rather than adding specialized layers or relying on extensive external resources, our method exploits *post-hoc decomposability* by examining and editing intermediate representations after the model is trained. By posing counterfactual questions (e.g., “What would the internal representation look like if the model aimed for a higher pitch rather than the current, lower one?”), we perform fine-grained modifications to hidden activations, thereby refining prosodic features and correcting mispronunciations at inference time.

The main contributions of this paper are summarized as follows: We introduce an approach that provides post-hoc adjustability of speech output, allowing for the adaptation of TTS models to new requirements without retraining. Additionally, we demonstrate the effectiveness of the proposed method through rigorous experimental validation.

2. Related Work

2.1. TTS Prosody Control

Text can be spoken in various ways due to semantic nuances, speaking styles, or inherent variability. Traditional approaches, such as unit-selection, capture this variability through speech databases [17]. In contrast, recent studies model prosodic variation by predicting key prosodic features such as pitch, duration, and energy from their embedding spaces [4, 5, 14, 18], or by using implicit representations learned from a reference encoder to capture the nuanced variations not specified by text alone [10, 11, 12].

Controlling intermediate representations has also been considered in prior work through the use of embedding bias, calculated by assessing the extent of translation required to achieve specific modifications in acoustic features within multidimensional scaling coordinates [19]. However, such modified representations risk deviating from the data manifold, and their application has been primarily confined to controlling duration.

2.2. TTS Pronunciation Control

A critical challenge for end-to-end TTS models, which operate without the aid of grapheme-to-phoneme dictionaries or predictive model, is polyphone disambiguation [20]. For TTS models to accurately convert graphemes into phonemes, their linguistic encoder must internalize the varied pronunciation rules, but

Audio samples can be found at <https://leekwoon.github.io/cae-tts>

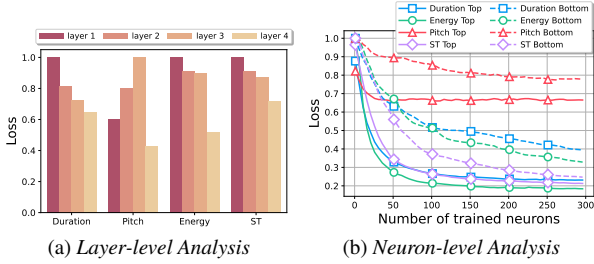


Figure 1: *Acoustic Correlation Analysis*. (a) *Classifier prediction losses across Tacotron 2 encoder layers, normalized with the highest value as one*. (b) *Effect of training only the most or least important neurons on classifier performance*.

fully internalizing them is difficult, leading to inevitable pronunciation errors in synthesized speech.

To mitigate mispronunciations, unit selection concatenates recorded speech fragments from a database. However, this often leads to noticeable join artifacts and requires a single-speaker database, limiting the use of non-target speaker data.

Recently, a model-centric approach called as the Speech Audio Corrector (SAC) [21] has been introduced. It leverages speech codes aligned with words, derived from self-supervised learning models to correct mispronunciations at the word level. In contrast, we introduce a model-agnostic approach that manipulates intermediate representations in TTS models.

3. Proposed Method

We first analyze intermediate representations produced by the encoder of an end-to-end TTS model, applicable to any encoder-decoder architectures.

Consider a sequence of input features $\mathbf{t} = \{t_1, \dots, t_n\}$, and let \mathbb{M} be a neural network model that transforms \mathbf{t} into a series of internal representations $\mathbf{t} \xrightarrow{\mathbb{M}} \mathbf{x} = \{x_1, \dots, x_n\}$. In this context, \mathbb{M} represents the encoder, with \mathbf{t} as either grapheme or phoneme embeddings. Let \mathbf{x}_i denote the activation in the last layer of encoder for the i -th input feature.¹

3.1. Acoustic Correlation Analysis

In the context of a classification task², we aim to predict a specific property, denoted as \mathbf{l}_i , from a set of properties \mathcal{P} presumed to be inherently learned by the model \mathbb{M} , for example, prosodic or pronunciation features in TTS model. Our objective is to examine intermediate representations within \mathbb{M} to assess their significance in recognizing the property $\mathbf{l}_i \in \mathcal{P}$.

For analyzing prosody properties such as duration, pitch, and energy, we first extract durations using an automatic speech recognition (ASR) system with connectionist temporal classification (CTC) alongside a backtranslation loss inspired by [22]. This system facilitates the extraction of precise alignments by sorting the posterior probabilities and applying monotonic alignment search (MAS). Subsequently, we calculate the average pitch and energy values based on these alignments to provide a comprehensive analysis.

For pronunciation analysis, we capture discrete pronunciation patterns at the semantic token (ST) level using the layer-6 representation of HuBERT [23], which operates at a 50Hz

¹While our analysis primarily concentrates on neurons in the top layer of the encoder, the approach can be applied to other layers.

²Prosody typically suits regression due to its continuous nature; however, for clarity, we discuss it under classification contexts here.

frame rate for 16kHz audio, to generate a sequence of semantic representations. These representations are aligned to input features using the same technique as for prosody, then aggregated and discretized into semantic tokens via k-means clustering (100 clusters). This process helps identify unique pronunciation patterns, supported by previous studies [21, 23, 24] that show HuBERT encodes phonetic information.

We train a logistic regression model on pairs of intermediate representations and their labels $\mathbf{x}_i, \mathbf{l}_i$, using cross-entropy loss. Previous studies have demonstrated that in the analysis of neural network representations, the performance characteristics of non-linear models closely resemble those of linear models [25, 26]. Our analysis reveals that TTS model representations are densely packed with information critical to both prosody and pronunciation, as shown in Figure 1.

3.2. Counterfactual Activation Editing

To manipulate the internal representations in a way that changes a predicted property, we consider a counterfactual question: for instance, “What would the internal representation look like if the synthesized speech had a higher pitch rather than a lower pitch?”. To address this, we introduce *Counterfactual Activation Editing (CAE)*, which manipulates internal representations to achieve the desired property.

Let $f : \mathcal{X} \rightarrow \mathbb{R}^C$ be a classifier that assigns to an activation vector $\mathbf{x}_i \in \mathcal{X}$ the probability of belonging to a class $c \in \{1, \dots, C\}$. For a given step size η and a target class c , CAE updates the activation via gradient ascent:

$$\mathbf{x}_i^{(k+1)} = \mathbf{x}_i^{(k)} + \eta \frac{\partial f_c}{\partial \mathbf{x}_i}(\mathbf{x}_i^{(k)}). \quad (1)$$

This process continues until the classifier output for the target class c exceeds a chosen threshold.³ The procedure can be formulated as:

$$\max_{\mathbf{x}'_i \in N(\mathbf{x}_i^{(k)}) \subset \mathbb{R}^d} f_c(\mathbf{x}'_i), \quad (2)$$

where $N(\mathbf{x}_i^{(k)}) = \{\mathbf{x} \in \mathbb{R}^d \mid d(\mathbf{x}, \mathbf{x}_i^{(k)}) < r\}$ denotes a neighborhood of radius r around the current activation $\mathbf{x}_i^{(k)}$, which depends on the optimization step size η .

3.2.1. Manifold Preserving CAE

In practice, the activation \mathbf{x}_i lies in a much lower-dimensional space, such as a k dimensional manifold \mathcal{M} , where $k \ll d$. The naive gradient-ascent procedure in the ambient space \mathbb{R}^d may move \mathbf{x}_i off the manifold, introducing unstructured noise rather than meaningful, semantically coherent changes. To mitigate this, we constrain updates to a local tangent space:

$$N_{\mathcal{T}}(\mathbf{x}_i^{(k)}) = \{\mathbf{x} \in \mathcal{T}_{\mathbf{x}_i^{(k)}} \mathcal{M} \mid d(\mathbf{x}, \mathbf{x}_i^{(k)}) < r\}, \quad (3)$$

where $\mathcal{T}_{\mathbf{x}_i^{(k)}} \mathcal{M}$ is the tangent space of \mathcal{M} at $\mathbf{x}_i^{(k)}$. We approximate these tangent spaces using an autoencoder, leveraging the information bottleneck to learn the manifold geometry [27, 28].

Concretely, we first map \mathbf{x}_i to the latent space \mathbf{z}_i of a β -VAE [29]. We then perform gradient ascent in \mathbf{z}_i , ensuring that the updates remain close to the data manifold. Formally,

$$\mathbf{z}_i^{(k+1)} = \mathbf{z}_i^{(k)} + \eta \frac{\partial (f \circ g)_c}{\partial \mathbf{z}_i}(\mathbf{z}_i^{(k)}), \quad (4)$$

³In a regression setting, where no explicit class boundaries exist, one may directly increase or decrease the regressor output toward the desired value.

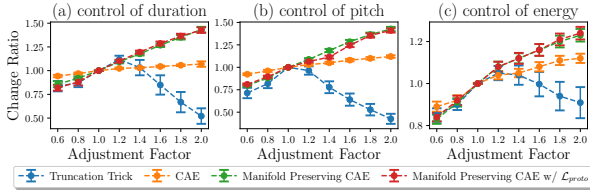


Figure 2: Ratio of change in prosodic features. Values above 1 indicate an increase, while values below 1 indicate a decrease.

where g is the decoder mapping from the latent space \mathcal{Z} back to the original activation space \mathcal{X} . By constraining edits within \mathcal{Z} , the resulting activations lie closer to the manifold \mathcal{M} , yielding counterfactual changes that are semantically meaningful.

3.2.2. Prototype Loss

Even with manifold preserving updates, the latent space of a β -VAE may include multiple subregions corresponding to different pronunciations or prosodic patterns. Consequently, large or unconstrained edits could inadvertently shift the activation into a subregion that alters the intended content. To address this, we incorporate a *prototype loss* that anchors the edited latent vectors to a representative code within a Vector Quantized VAE (VQ-VAE) [30].

Let $\text{Enc}(\cdot)$ and $\text{Dec}(\cdot)$ be the encoder and decoder of the VQ-VAE, and let $\mathbf{e} \in \mathbb{R}^{K \times D}$ denote the learned codebook with K discrete embeddings of dimension D . Given a latent representation \mathbf{z}_i obtained via our Manifold Preserving CAE process, we find its nearest codebook vector \mathbf{e}_k :

$$\begin{aligned} \mathbf{e}_k &= \text{Quantize}(\text{Enc}(\mathbf{z}_i)) \\ &= \arg \min_j \|\text{Enc}(\mathbf{z}_i) - \mathbf{e}_j\|_2. \end{aligned} \quad (5)$$

Decoding \mathbf{e}_k yields the *prototype* representation:

$$\text{proto}(\mathbf{z}_i) = \text{Dec}(\mathbf{e}_k) \quad (6)$$

To ensure edits remain faithful to the intended pronunciation, we augment the existing objective with a prototype loss that encourages \mathbf{z}_i to remain close to its prototype:

$$\mathcal{L}_{\text{proto}} = \alpha \cdot \|\mathbf{z}_i - \text{proto}(\mathbf{z}_i)\|_2^2. \quad (7)$$

where α is a weighting factor. By enforcing proximity to a prototype in the codebook, this loss discourages transitions to latent subregions that would alter the target speech content.

4. Experimental Setup

4.1. TTS Model

We experimented with Tacotron 2 [2] encoder-decoder architecture, a model with about 28 million parameters, including an encoder with three convolutional layers followed by an LSTM layer. Although Tacotron 2 generates highly natural speech, it lacks explicit mechanisms to capture diverse prosodic contours without additional modules or training data. This limitation makes Tacotron 2 an ideal testbed for our *post-hoc* modification approach, as any improvement in controllability can be primarily attributed to our method rather than inherent model capacity. Phoneme inputs are employed to control prosody, whereas grapheme inputs are used for pronunciation adjustments, to better represent resource-limited scenarios. We base our Tacotron

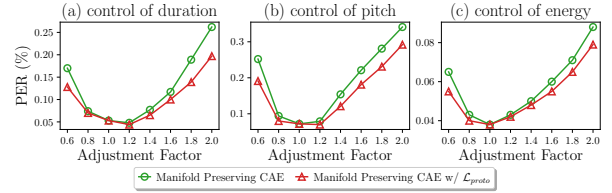


Figure 3: Phone-Error-Rate (PER) comparison with and without prototype loss, measured using the Whisper ASR system.

2 on NVIDIA’s implementation⁴ and use a pre-trained WaveGlow [7] as the vocoder.

4.2. Dataset

For training Tacotron 2, we use the LJSpeech dataset [31], containing 13,100 English audio clips (around 24 hours) from a single female speaker, along with text transcripts. The dataset was split into 12,500 samples for training, 500 for validation, and 100 for testing.

4.3. Training

In the training phase of the classifier, we extracted acoustic properties and semantic tokens as detailed in Section 3.1 from the LJSpeech training dataset. For prosody prediction, a linear regression model was employed, using mean squared error loss augmented with elastic net regularization, setting the l1 and l2 regularization hyperparameters, λ_1 and λ_2 , both to 0.001. For semantic token prediction, we used a logistic regression model trained using cross-entropy loss. Both models were trained using a learning rate of 0.01. Additionally, for the training of the β -VAE, we set β to 0.04 and the dimension of the latent space to 16, with a learning rate of 0.0002.

5. Results

Our experiments study the following questions: (1) Are acoustic features encoded by the Tacotron 2 encoder? (2) Can *counterfactual activation editing* (CAE) effectively manipulate prosody? (3) Can CAE further correct mispronunciations?

5.1. Acoustic Correlation Analysis

To assess the capacity of Tacotron 2 encoder layers to encapsulate meaningful acoustic features, classifiers were trained using activations from different layers of a pre-trained Tacotron 2 model. Prediction losses across the encoder layers, shown in Figure 1a, demonstrate that acoustic information is captured at varying levels. Notably, the last LSTM layer is more accurate in capturing acoustic properties compared to the earlier convolutional layers.

In addition, we conduct a detailed examination of individual neurons to identify those that significantly influence certain properties. This involved training only the most or least significant neurons from the last layer of the encoder, ranked by their importance based on classifier weights [26], and observing the effect on classifier efficacy. Figure 1b shows that classifier performance improves more significantly when top-ranked neurons are trained (solid lines) compared to when the least important neurons are trained (dashed lines), indicating the crucial role of specific neurons in encoding acoustic properties.

⁴<https://github.com/NVIDIA/tacotron2>

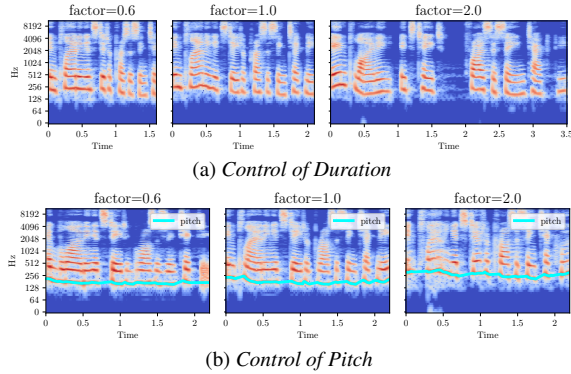


Figure 4: Spectrograms of the synthesized speech using Manifold Preserving CAE for prosody control. The input text used for synthesis is “In planning its data processing techniques.”

5.2. Prosody Control

We investigate *Counterfactual Activation Editing* (CAE) for controlling prosodic features in a pre-trained Tacotron 2 model. We compare three variations: CAE, which edits the encoder activations directly; Manifold Preserving CAE, which first encodes activations into a VAE latent space and edits them there to prevent manifold deviation; and Manifold Preserving CAE with prototype loss, which further constrains edits around representative codebook vectors to preserve intelligibility. We also compare against a truncation trick [32], which shifts activations toward values observed in the training set.

We specifically modify the final-layer encoder activations for each phoneme in order to scale its prosodic property (e.g., pitch, duration, or energy) by a predefined adjustment factor. Let $f(\mathbf{x}_i)$ be a regressor predicting a prosodic feature from the activation vector \mathbf{x}_i . Suppose $f(\mathbf{x}_i) = p_i$ is the current predicted pitch. To increase or decrease pitch by a factor $\lambda > 0$, we adjust \mathbf{x}_i so that $f(\mathbf{x}'_i)$ approaches λp_i , following the gradient-based procedure described in Section 3.2.

Figure 2 shows the ratio between the modified prosodic feature and its original value, where values above 1 indicate an increase and values below 1 indicate a decrease. Manifold Preserving CAE consistently achieves more precise prosodic changes compared to both CAE and the truncation trick. In particular, the truncation trick often led to omissions or mispronunciations. Figure 4 illustrates how the spectrogram changes when pitch or duration is controlled. We used Whisper [33] to measure the Phone-Error-Rate (PER), finding that prototype loss plays a crucial role in maintaining speech intelligibility under large prosodic shifts (see Figure 3). Additionally, the edits can be localized to specific phonemes or words by selecting only the corresponding encoder activations for adjustment, as shown in the audio examples on our demo webpage.

5.3. Correction of Mispronunciations

To correct mispronunciations within synthesized speech, we assume supervision in the form of a speech-only correction query. The premise is that word-level pronunciation information, more easily obtainable than phonetic details, can serve as a practical foundation for corrective adjustments.

Building on the prosody control method, we target mispronunciation correction by manipulating activations in the final Tacotron 2 encoder layer. Our goal is to adjust the predicted semantic tokens for each grapheme to match those of the correction query, ensuring more accurate pronunciation.

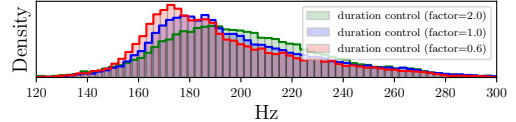


Figure 5: Phoneme pitch density over the LJSpeech dataset after Manifold Preserving CAE duration control.

To evaluate pronunciation correction, we synthesized 78 speech samples, each with a target word in the carrier sentence “How is ... pronounced?”. These target words were selected from prior work [21] as being outside the LJSpeech training set or challenging for grapheme-input TTS models.

Mispronunciation correction, using Manifold Preserving CAE, was evaluated using four metrics: Word-Error-Rate (WER), Phone-Error-Rate (PER), sentence embedding cosine similarity, and token-based similarity. For cosine similarity, we converted synthesized speech to text using Whisper ASR, then compared sentence embeddings from the text-ada-002 model to ground truth. Token-based similarity was measured using the GPT-2 [34] tokenizer on Whisper transcripts, assessing the ratio of correctly identified tokens to the ground truth.

Table 1: The objective correction capability comparisons.

Metric	Before Correction	After Correction
WER (\downarrow)	0.151	0.056
PER (\downarrow)	0.069	0.017
Cosine Similarity (\uparrow)	0.954	0.969
Token Similarity (\uparrow)	0.892	0.929

As indicated in Table 1, our CAE approach is effective across all pronunciation correction metrics, substantially reducing both WER and PER, while increasing the semantic similarity of the corrected outputs.

For evaluating the perceptual quality, we conducted a Comparison Mean Opinion Score (CMOS) [35] test comparing initial synthesized speech with corrected speech. Twenty listeners rated the quality, with results presented in Table 2.

Table 2: CMOS comparison.

Setting	CMOS
After Correction	0.000
Before Correction	-0.764

The results show that correcting prosody and mispronunciations with our CAE approach led to a 0.764-point improvement in CMOS, demonstrating its effectiveness.

6. Conclusions and Discussion

In this work, we introduced a novel method for prosody control and pronunciation correction in TTS models using Counterfactual Activation Editing, enabling post-hoc adjustments without retraining. Our experiments demonstrated that intermediate representations within TTS models, specifically Tacotron 2, contain rich information that can be effectively manipulated to modify prosody and correct mispronunciations.

Despite its effectiveness, our approach manipulates entire vectors of intermediate representations, which can overlook individual dimension contributions. For instance, as shown in Figure 5, controlling duration slightly affects pitch, suggesting feature entanglement. Future investigations will aim to refine our approach by targeting specific dimensions (neurons) within these representations. This direction holds the promise of achieving a more granular, interpretable, and disentangled control over speech characteristics.

7. Acknowledgements

This work was supported by Samsung Electronics MX Division and by Institute of Information & communications Technology Planning & Evaluation (IITP) grant funded by the Korea government (MSIT) (No.2019-0-00075, Artificial Intelligence Graduate School Program (KAIST); No. 2022-0-00984, Development of Artificial Intelligence Technology for Personalized Plug-and-Play Explanation and Verification of Explanation; No. RS-2024-00457882, AI Research Hub Project).

8. References

- [1] 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.
- [2] 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*. IEEE, 2018, pp. 4779–4783.
- [3] N. Li, S. Liu, Y. Liu, S. Zhao, and M. Liu, “Neural speech synthesis with transformer network,” in *Proceedings of the AAAI conference on artificial intelligence*, vol. 33, no. 01, 2019, pp. 6706–6713.
- [4] 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.
- [5] 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,” *arXiv preprint arXiv:2006.04558*, 2020.
- [6] 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.
- [7] R. Prenger, R. Valle, and B. Catanzaro, “Waveglow: A flow-based generative network for speech synthesis,” in *ICASSP*. IEEE, 2019, pp. 3617–3621.
- [8] K. Kumar, R. Kumar, T. De Boissiere, L. Gestin, W. Z. Teoh, J. Sotelo, A. De Brebisson, Y. Bengio, and A. C. Courville, “Melgan: Generative adversarial networks for conditional waveform synthesis,” *Advances in neural information processing systems*, vol. 32, 2019.
- [9] J. Kong, J. Kim, and J. Bae, “Hifi-gan: Generative adversarial networks for efficient and high fidelity speech synthesis,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 17 022–17 033, 2020.
- [10] Y. Wang, D. Stanton, Y. Zhang, R.-S. Ryan, E. Battenberg, J. Shor, Y. Xiao, Y. Jia, F. Ren, and R. A. Saurous, “Style tokens: Unsupervised style modeling, control and transfer in end-to-end speech synthesis,” in *International conference on machine learning*. PMLR, 2018, pp. 5180–5189.
- [11] R. Skerry-Ryan, E. Battenberg, Y. Xiao, Y. Wang, D. Stanton, J. Shor, R. 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*. PMLR, 2018, pp. 4693–4702.
- [12] W.-N. Hsu, Y. Zhang, R. J. Weiss, H. Zen, Y. Wu, Y. Wang, Y. Cao, Y. Jia, Z. Chen, J. Shen *et al.*, “Hierarchical generative modeling for controllable speech synthesis,” *arXiv preprint arXiv:1810.07217*, 2018.
- [13] J. Šimko, T. Törö, M. Vainio, and A. Suni, “Prosody under control: Controlling prosody in text-to-speech synthesis by adjustments in latent reference space,” in *Proceedings of the 20th International Congress of Phonetic Sciences*, 2023, pp. 3086–3090.
- [14] D. S. R. Mohan, V. Hu, T. H. Teh, A. Torresquintero, C. G. Wallis, M. Staib, L. Foglianti, J. Gao, and S. King, “Ctrl-p: Temporal control of prosodic variation for speech synthesis,” *arXiv preprint arXiv:2106.08352*, 2021.
- [15] T. Raitio, J. Li, and S. Seshadri, “Hierarchical prosody modeling and control in non-autoregressive parallel neural tts,” in *ICASSP*. IEEE, 2022, pp. 7587–7591.
- [16] J. Taylor and K. Richmond, “Analysis of pronunciation learning in end-to-end speech synthesis,” in *Interspeech*, 2019.
- [17] V. Strom, R. A. Clark, and S. King, “Expressive prosody for unit-selection speech synthesis,” in *Interspeech*, 2006.
- [18] J. Bandekar, S. Udupa, A. Singh, A. Jayakumar, S. Badiger, S. Kumar, P. VH, P. K. Ghosh *et al.*, “Speaking rate attention-based duration prediction for speed control tts,” *arXiv preprint arXiv:2310.08846*, 2023.
- [19] M. Lenglet, O. Perrotin, and G. Bailly, “Speaking rate control of end-to-end tts models by direct manipulation of the encoder’s output embeddings,” in *Interspeech*, 2022.
- [20] Y. Zhang, L. Deng, and Y. Wang, “Unified mandarin tts front-end based on distilled bert model,” *arXiv preprint arXiv:2012.15404*, 2020.
- [21] J. Fong, D. Lyth, G. E. Henter, H. Tang, and S. King, “Speech audio corrector: using speech from non-target speakers for one-off correction of mispronunciations in grapheme-input text-to-speech,” in *Interspeech*, 2022.
- [22] F. Lux, J. Koch, and N. T. Vu, “Exact prosody cloning in zero-shot multispeaker text-to-speech,” in *2022 IEEE Spoken Language Technology Workshop (SLT)*. IEEE, 2023, pp. 962–969.
- [23] W.-N. Hsu, B. Bolte, Y.-H. H. Tsai, K. Lakhota, R. Salakhutdinov, and A. Mohamed, “Hubert: Self-supervised speech representation learning by masked prediction of hidden units,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 3451–3460, 2021.
- [24] H. Ji, T. Patel, and O. Scharenborg, “Predicting within and across language phoneme recognition performance of self-supervised learning speech pre-trained models,” *arXiv preprint arXiv:2206.12489*, 2022.
- [25] Y. Belinkov, N. Durrani, F. Dalvi, H. Sajjad, and J. Glass, “What do neural machine translation models learn about morphology?” *arXiv preprint arXiv:1704.03471*, 2017.
- [26] F. Dalvi, N. Durrani, H. Sajjad, Y. Belinkov, A. Bau, and J. Glass, “What is one grain of sand in the desert? analyzing individual neurons in deep nlp models,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, 2019, pp. 6309–6317.
- [27] S. Joshi, O. Koyejo, W. Vijitbenjaronk, B. Kim, and J. Ghosh, “Towards realistic individual recourse and actionable explanations in black-box decision making systems,” *arXiv preprint arXiv:1907.09615*, 2019.
- [28] A.-K. Dombrowski, J. E. Gerken, K.-R. Müller, and P. Kessel, “Diffeomorphic counterfactuals with generative models,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- [29] I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, and A. Lerchner, “beta-vae: Learning basic visual concepts with a constrained variational framework,” in *International conference on learning representations*, 2016.
- [30] A. Van Den Oord, O. Vinyals *et al.*, “Neural discrete representation learning,” *Advances in neural information processing systems*, vol. 30, 2017.
- [31] K. Ito and L. Johnson, “The lj speech dataset,” <https://keithito.com/LJ-Speech-Dataset/>, 2017.
- [32] A. Brock, J. Donahue, and K. Simonyan, “Large scale gan training for high fidelity natural image synthesis,” *arXiv preprint arXiv:1809.11096*, 2018.
- [33] A. Radford, J. W. Kim, T. Xu, G. Brockman, C. McLeavey, and I. Sutskever, “Robust speech recognition via large-scale weak supervision,” in *International Conference on Machine Learning*. PMLR, 2023, pp. 28 492–28 518.
- [34] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever *et al.*, “Language models are unsupervised multitask learners,” *OpenAI blog*, vol. 1, no. 8, p. 9, 2019.
- [35] P. C. Loizou, “Speech quality assessment,” in *Multimedia analysis, processing and communications*. Springer, 2011, pp. 623–654.