# Towards Controllable Speech Synthesis in the Era of Large Language Models: A Survey

Tianxin Xie, Yan Rong, Pengfei Zhang, Li Liu

*Abstract*—Text-to-speech (TTS), also known as speech synthesis, is a prominent research area that aims to generate natural-sounding human speech from text. Recently, with the increasing industrial demand, TTS technologies have evolved beyond synthesizing human-like speech to enabling controllable speech generation. This includes fine-grained control over various attributes of synthesized speech such as emotion, prosody, timbre, and duration. Besides, advancements in deep learning, such as diffusion and large language models, have significantly enhanced controllable TTS over the past several years. In this paper, we conduct a comprehensive survey of controllable TTS, covering approaches ranging from basic control techniques to methods utilizing natural language prompts, aiming to provide a clear understanding of the current state of research. We examine the general controllable TTS pipeline, challenges, model architectures, and control strategies, offering a comprehensive and clear taxonomy of existing methods. Additionally, we provide a detailed summary of datasets and evaluation metrics and shed some light on the applications and future directions of controllable TTS. To the best of our knowledge, this survey paper provides the first comprehensive review of emerging controllable TTS methods, which can serve as a beneficial resource for both academic researchers and industry practitioners.

*Index Terms*—Text-to-speech, controllable TTS, speech synthesis, TTS survey, large language models, diffusion models.

## I. INTRODUCTION

Speech synthesis, also broadly known as text-to-speech (TTS), is a long-time developed technique that aims to synthesize human-like voices from text [1], [2], and it has extensive applications in our daily lives, such as health care [3], [4], personal assistants [5], entertainment [6], [7], and robotics [8], [9]. Recently, TTS has gained significant attention with the rise of large language model (LLM)-powered chatbots, such as ChatGPT [10] and Llama [11], due to its naturalness and convenience for human-computer interaction. Meanwhile, the ability to achieve fine-grained control over synthesized speech attributes, such as emotion, prosody, timbre, and duration, has become a hot research topic in both academia and industry, driven by its vast potential for diverse applications.

Deep learning [12] has made great progress in the past decade due to exponentially growing computational resources like GPUs [13], leading to the explosion of numerous great works on TTS [14]–[17]. These methods can synthesize human speech with better quality [14] and can achieve fine-grained control of the generated voice [18]–[22]. Besides, some recent

works synthesize speech given multi-modal input, such as face images [23], [24], cartoons [7], and videos [25]. Moreover, with the fast development of open-source LLMs [11], [26]–[29], some researchers propose to synthesize fine-grained controllable speech with natural language description [30]–[32], coining a new way to generate custom speech voices. Meanwhile, powering LLMs with speech synthesis has also been a hot topic in the last few years [33]–[35]. In recent years, a wide range of TTS methods has emerged, making it essential for researchers to gain a comprehensive understanding of current research trends, particularly in controllable TTS, to identify promising future directions in this rapidly evolving field. Consequently, there is a pressing need for an up-to-date survey of TTS techniques. While several existing surveys address parametric-based approaches [36]–[41] and deep learning-based TTS [42]–[48], they largely overlook the controllability of TTS. Additionally, these surveys do not cover the advancements in recent years, such as natural language description-based TTS methods.

This paper provides a comprehensive and in-depth survey of existing and emerging TTS technologies, with a particular focus on controllable TTS methods. Fig. 1 demonstrates the development of controllable TTS methods in recent years, showing their backbones, feature representations, and control abilities. The remainder of this section begins with a brief comparison between this survey and previous ones, followed by an overview of the history of controllable TTS technologies, ranging from early milestones to state-of-the-art advancements. Finally, we introduce the taxonomy and organization of this paper.

### A. Comparison with Existing Surveys

Several survey papers have reviewed TTS technologies, spanning early approaches from previous decades [36], [37], [40], [49] to more recent advancements [42], [43], [50]. However, to the best of our knowledge, this paper is the first to focus specifically on controllable TTS. The key differences between this survey and prior work are summarized as follows:

**Different scope.** Klatt et al. [36] provided the first comprehensive survey on formant, concatenative, and articulatory TTS methods, with a strong emphasis on text analysis. In the early 2010s, Tabet et al. [49] and King et al. [40] explored rule-based, concatenative, and HMM-based techniques. Later, the advent of deep learning catalyzed the emergence of numerous neural-based TTS methods. Therefore, Ning et al. [43] and Tan et al. [42] have conducted extensive surveys on neural-based acoustic models and vocoders, while Zhang
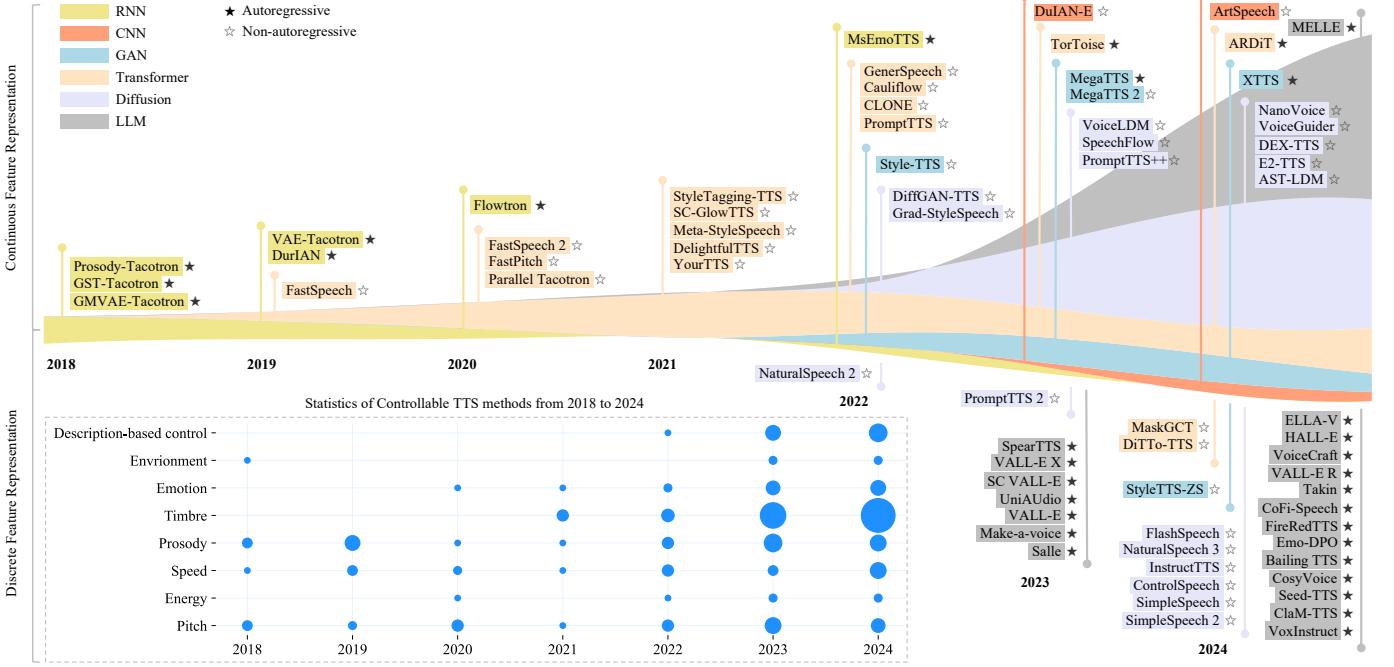
Fig. 1. A summary of representative controllable TTS methods in recent years and their model architectures, feature representations, and control abilities. There are additional network structures, such as VAE-based models, that are not included in this figure. For more details, refer to Tables IV and III.

et al. [50] presented the first review of diffusion model-based TTS techniques. However, these studies offer limited discussion on the controllability of TTS systems. To address this gap, we present the first comprehensive survey of TTS methods through the lens of controllability, providing an in-depth analysis of model architectures and strategies for controlling synthesized speech.

**Close to current demand.** With the rapid development of hardware (i.e., GPUs) and artificial intelligence techniques (i.e., transformers, LLMs, diffusion models) in the last few years, the demand for controllable TTS is becoming increasingly urgent due to its broad applications in industries such as filmmaking, gaming, robots, and personal assistants. Despite this growing need, existing surveys pay little attention to control methods in TTS technologies. To bridge this gap, we propose a systematic analysis of current controllable TTS methods and the associated challenges, offering a comprehensive understanding of the research state in this field.

**New insights and directions.** This survey offers new insights through a comprehensive analysis of model architectures and control methods in controllable TTS systems. Additionally, it provides an in-depth discussion of the challenges associated with various controllable TTS tasks. Furthermore, we address the question: "Where are we on the path to fully controllable TTS technologies?", by examining the relationship and gap between current TTS methods and industrial requirements. Based on these analyses, we identify promising directions for future research on TTS technologies.

Table I summarizes representative surveys and this paper in terms of main focus and publication year.

### TABLE I
### COMPARISON WITH REPRESENTATIVE TTS SURVEYS.

| Survey | Main Focus | Year |
|---|---|---|
| Klatt et al. [36] | Rule-based and concatenative TTS | 1987 |
| Tabet et al. [49] | Rule-based, concatenative, and parametric TTS | 2011 |
| King et al. [40] | Parametric TTS and performance measurement | 2014 |
| Tan et al. [42] | Neural-based, efficient, and expressive TTS | 2021 |
| Zhang et al. [50] | Diffusion-based TTS and speech enhancement | 2023 |
| Ours | Controllable TTS and evaluation | 2024 |

### B. The History of Controllable TTS

Controllable TTS aims to control various aspects of synthesized speech, such as pitch, energy, speed/duration, prosody, timbre, emotion, gender, or high-level styles. This subsection briefly reviews the history of controllable TTS ranging from early approaches to the state-of-arts in recent years.

**Early approaches.** Before the prevalence of deep neural networks (DNNs), controllable TTS technologies were built primarily on rule-based, concatenative, and statistical methods. These approaches enable some degree of customization and control, though they were constrained by the limitations of the underlying models and available computational resources. 1) Rule-based TTS systems [51]–[54], such as formant synthesis, were among the earliest methods for speech generation. These systems use manually crafted rules to simulate the speech generation process by controlling acoustic parameters such as pitch, duration, and formant frequencies, allowing explicit manipulation of prosody and phonetic details through rule adjustments. 2) Concatenative TTS [55]–[58], which dominated the field in the late 1990s and early 2000s, synthesize speech by concatenating pre-recorded speech segments, such

as phonemes or diphones, stored in a large database [59]. These methods can modify the prosody by manipulating the pitch, duration, and amplitude of speech segments during concatenation. They also allow limited voice customization by selecting speech units from different speakers. 3) Parametric methods, particularly HMM-based TTS [60]–[65], gained prominence in the late 2000s. These systems model the relationships between linguistic features and acoustic parameters, providing more flexibility in controlling prosody, pitch, speaking rate, and timbre by adjusting statistical parameters. Some HMM-based systems also supported speaker adaptation [66], [67] and voice conversion [68], [69], enabling voice cloning to some extent. Besides, emotion can also be limitedly controlled by some of these methods [60], [70]–[72]. In addition, they required less storage compared to concatenative TTS and allowed smoother transitions between speech units.

**Neural-based synthesis.** Neural-based TTS technologies emerged with the advent of deep learning, significantly advancing the field by enabling more flexible, natural, and expressive speech synthesis. Unlike traditional methods, neural-based TTS leverages DNNs to model complex relationships between input text and speech, facilitating nuanced control over various speech characteristics. Early neural TTS systems, such as WaveNet [73] and Tacotron [74] laid the groundwork for controllability. 1) Controlling prosody features like rhythm and intonation is vital for generating expressive and contextually appropriate speech. Neural-based TTS models achieve prosody control through explicit conditioning or learned latent representations [15], [75]–[78]. 2) Speaker control has also gained significant improvement in neural-based TTS through speaker embeddings or adaptation techniques [79]–[82]. 3) Besides, emotionally controllable TTS [20], [22], [31], [32], [83] has become a hot topic due to the strong modeling capability of DNNs, enabling the synthesis of speech with specific emotional tones such as happiness, sadness, anger, or neutrality. These systems go beyond producing intelligible and natural-sounding speech, focusing on generating expressive output that aligns with the intended emotional context. 4) Neural-based TTS can also manipulate timbre (vocal quality) [14], [78], [84]–[87] and style (speech mannerisms) [88]–[90], allowing for creative and personalized applications. These techniques lead to one of the most popular research topics, i.e., zero-shot TTS (particularly voice cloning) [78], [82], [91], [92]. 5) Fine-grained content and linguistic control also become more powerful [93]–[96]. These methods can emphasize or de-emphasize specific words or adjust the pronunciation of phonemes through speech editing or generation techniques.

Neural-based TTS technologies represent a significant leap in the flexibility and quality of speech synthesis. From prosody and emotion to speaker identity and style, these systems empower diverse applications in fields such as entertainment, accessibility, and human-computer interaction.

**LLM-based synthesis.** Here we pay special attention to LLM-based synthesis methods due to their superior context modeling capabilities compared to other neural-based TTS methods. LLMs, such as GPT [97], [98], T5 [99], and PaLM [100], have revolutionized various natural language processing (NLP) tasks with their ability to generate coherent, context-aware text. Recently, their utility has expanded into controllable TTS technologies [17], [101]–[104]. For example, users can synthesize the target speech by describing its characteristics, such as: "A young girl says 'I really like it, thank you!' with a happy voice", making speech generation significantly more intuitive and user-friendly. Specifically, an LLM can detect emotional intent in sentences (e.g., "I'm thrilled" → happiness, "This is unfortunate" → sadness). The detected emotion is encoded as an auxiliary input to the TTS model, enabling modulation of acoustic features like prosody, pitch, and energy to align with the expressed sentiment. By leveraging LLMs' capabilities in understanding and generating rich contextual information, these systems can achieve enhanced and fine-grained control over various speech attributes such as prosody, emotion, style, and speaker characteristics [31], [105], [106]. Integrating LLMs into TTS systems represents a significant step forward, enabling more dynamic and expressive speech synthesis.

### C. Organization of This Survey

This paper first presents a comprehensive and systematic review of controllable TTS technologies, with a particular focus on model architectures, control methodologies, and feature representations. To establish a foundational understanding, this survey begins with an introduction to the TTS pipeline in Section II. While our focus remains on controllable TTS, Section III examines seminal works in uncontrollable TTS that have significantly influenced the field's development. Section IV provides a thorough investigation into controllable TTS methods, analyzing both their model architectures and control strategies. Section V presents a comprehensive review of datasets and evaluation metrics. Section VI provides an in-depth analysis of the challenges encountered in achieving controllable TTS systems and discusses future directions. Section VII explores the broader impacts of controllable TTS technologies and identifies promising future research directions, followed by the conclusion in Section VIII.

## II. TTS PIPELINE

In this section, we elaborate on the general pipeline that supports controllable TTS technologies, including acoustic models, speech vocoders, and feature representations. Fig.2 depicts the general pipeline of controllable TTS, containing various model architectures and feature representations, but the control strategies will be discussed in Section IV. Readers can jump to Section III if familiar with TTS pipelines.

### A. Overview

A TTS pipeline generally contains three key components, i.e., linguistic analyzer, acoustic model, speech vocoder, and with a conditional input, e.g., prompts, for controllable speech synthesis. Besides, some end-to-end methods use a single model to encode the input and decode the speech waveforms without generating intermediate features like mel-spectrograms [110]. *Linguistic analyzer* aims to extract linguistic features, e.g., phoneme duration and position, syllable
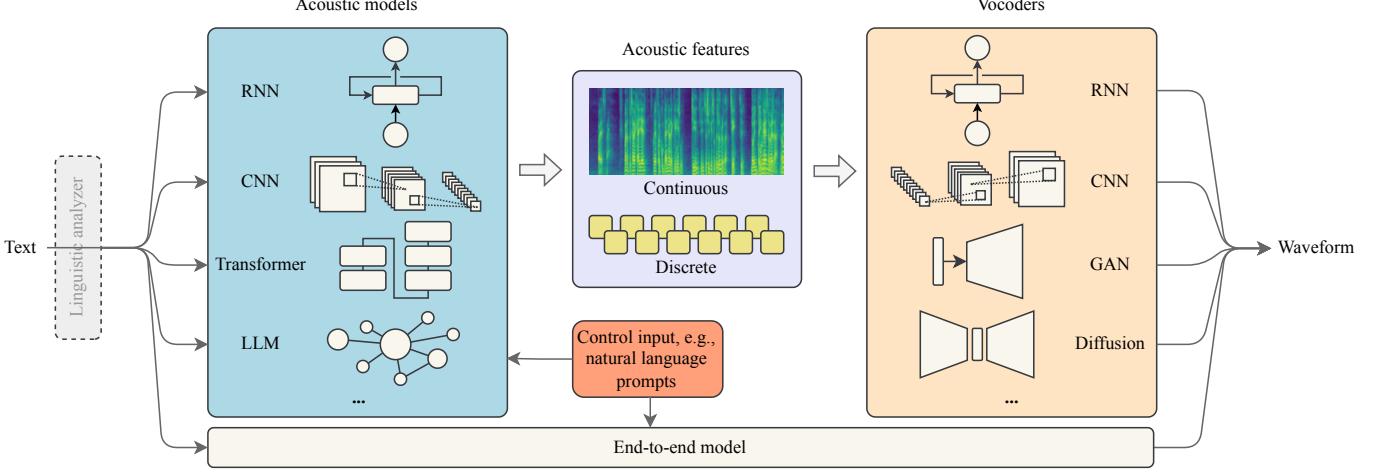
Fig. 2. The general pipeline of controllable TTS from the perspective of network structure. Linguistic analysis is necessary for parametric and a few neural-based methods but is no longer needed for most modern neural-based methods. In this paper, we only review neural-based controllable TTS methods and do not investigate acoustic features (e.g., MCC [107], LSP [108], F0 [109]) used in early TTS methods.

stress, and utterance level, from the input text, which is a necessary step in HHM-based methods [64], [65] and a few neural-based methods [111], [112], but is time-consuming and error-prone. *Acoustic model* is a parametric or neural model that predicts the acoustic features from the input texts. Modern neural-based acoustic models like Tacotron [74] and later works [15], [76], [113] directly take character [114] or word embeddings [115] as the input, which is much more efficient than previous methods. *Speech vocoder* is the last component that converts the intermediate acoustic features into a waveform that can be played back. This step bridges the gap between the acoustic features and the actual sounds produced, helping to generate high-quality, natural-sounding speech [73], [116]. Tan et al. [42] have presented a comprehensive and detailed review of acoustic models and vocoders. Therefore, the following subsections will briefly introduce some representative acoustic models and speech vocoders, followed by a discussion of acoustic feature representations.

### B. Acoustic Models

Acoustic modeling is a crucial step in TTS because it ensures the generated acoustic features capture the subtleties of human speech. By accurately modeling acoustic features, modern TTS systems can help generate high-quality and expressive audio that sounds close to human speech.

**Parametric models.** Early acoustic models rely on parametric approaches, where predefined rules and mathematical functions are utilized to model speech generation. These models often utilize HMMs to capture acoustic features from linguistic input and generate acoustic features by parameterizing the vocal tract and its physiological properties such as pitch and prosody [71], [72], [117]–[120]. These methods have relatively low computational costs and can produce a range of voices by adjusting model parameters. However, the speech quality of these methods is robotic and lacks natural intonation, and the expressiveness is also limited [72], [120].

**RNN-based models.** Recurrent Neural Networks (RNNs) proved particularly effective in early neural-based TTS due to their ability to model sequential data and long-range dependencies, which helps in capturing the sequential nature of speech, such as the duration and natural flow of phonemes. Typically, these models have an encoder-decoder architecture, where an encoder encodes input linguistic features, such as phonemes or text, into a fixed-dimensional representation, and the decoder sequentially decodes this representation into acoustic features (e.g., mel-spectrogram frames) that capture the frequency and amplitude of sound over time. Tacotron 2 [75] is one of the pioneering TTS models that uses RNNs with an attention mechanism, which helps align the text sequence with the generated acoustic features. It takes raw characters as input and produces mel-spectrogram frames, which are subsequently converted to waveforms. Another example is MelNet [121], which leverages autoregressive modeling to generate high-quality mel-spectrograms, demonstrating versatility in generating both speech and music, achieving high fidelity and coherence across temporal scales.

**CNN-based models.** Unlike RNNs, which process sequential data frame by frame, CNNs process the entire sequence at once by applying filters across the input texts. This parallel approach enables faster training and inference, making CNN-based TTS particularly appealing for real-time and low-latency applications. Furthermore, by stacking multiple convolutional layers with varying kernel sizes or dilation rates, CNNs can capture both short-range and long-range dependencies, which are essential for natural-sounding speech synthesis. Deep Voice [122] is one of the first prominent CNN-based TTS models by Baidu, designed to generate mel-spectrograms directly from phoneme or character input. ParaNet [123] also utilizes a RNN model to achieve sequence-to-sequence mel-spectrogram generation. It uses a non-autoregressive architecture, which enables significantly faster inference by predicting multiple time steps simultaneously.

**Transformer-based models.** Transformer model [124] uses

self-attention layers to capture relationships within the input sequence, making them well-suited for tasks requiring an understanding of global contexts, such as prosody and rhythm in TTS. Transformer-based TTS models often employ an encoder-decoder architecture, where the encoder processes linguistic information (e.g., phonemes or text) and captures contextual relationships, and the decoder generates acoustic features (like mel-spectrograms) from these encoded representations, later converted to waveforms by a vocoder. TransformerTTS [125] is one of the first TTS models that apply transformers to synthesize speech from text. It utilizes a standard encoder-decoder transformer architecture and relies on multi-head self-attention mechanisms to model long-term dependencies, which helps maintain consistency and natural flow in speech over long utterances. FastSpeech [15] is a non-autoregressive model designed to overcome the limitations of autoregressive transformers in TTS, achieving faster synthesis than previous methods. It introduces a length regulator to align text with output frames, enabling the control of phoneme duration. FastSpeech 2 [76] extends FastSpeech by adding pitch, duration, and energy predictors, resulting in more expressive and natural-sounding speech.

**LLM-based models.** LLMs [11], [26], [97], [126], known for their large-scale pre-training on text data, have shown remarkable capabilities in natural language understanding and generation. LLM-based TTS models generally use a text description to guide the mel-spectrogram generation, where the acoustic model processes the input text to generate acoustic tokens that capture linguistic and contextual information, such as tone, sentiment, and prosody. For example, PromptTTS [101] uses a textual prompt encoded by BERT [126] to guide the acoustic model on the timbre, tone, emotion, and prosody desired in the speech output. PromptTTS first generates mel-spectrograms with token embeddings and then converts them to audio using a vocoder. InstructTTS [105] generates expressive and controllable speech using natural language style prompts. It leverages discrete latent representations of speech and integrates natural language descriptions to guide the synthesis process, which bridges the gap between TTS systems and natural language interfaces, enabling fine-grained style control through intuitive prompts.

**Other acoustic models.** In TTS, GANs [127]–[129], VAEs [18], [130], and diffusion models [113], [131] can also be used as acoustic models. Flow-based methods [132], [133] are also popular in waveform generation. Refer to the survey paper from Tan et al. [42] for more details.

The choice of an acoustic model depends on the specific requirements and is a trade-off between synthesis quality, computational efficiency, and flexibility. For real-time applications, CNN-based or lightweight transformer-based models are preferable, while for high-fidelity, expressive speech synthesis, transformer-based and LLM-based models are better suited.

### C. Speech Vocoders

Vocoders are essential for converting acoustic features such as mel-spectrograms into intelligible audio waveforms and are vital in determining the naturalness and quality of synthesized speech. We broadly categorize existing vocoders according to their model architectures, i.e., RNN-, CNN-, GAN-, and diffusion-based vocoders.

**RNN-based vocoders.** Unlike traditional vocoders [134], [135] that depend on manually designed signal processing pipelines, RNN-based vocoders [136]–[139] leverage the temporal modeling capabilities of RNNs to directly learn the complex patterns in speech signals, enabling the synthesis of natural-sounding waveforms with improved prosody and temporal coherence. For instance, WaveRNN [137] generates speech waveforms sample-by-sample using a single-layer recurrent neural network, typically with Gated Recurrent Units (GRU). It improves upon earlier neural vocoders like WaveNet [73] by significantly reducing the computational requirements without sacrificing audio quality. MB-WaveRNN [139] extends WaveRNN by incorporating a multi-band decomposition strategy, where the speech waveform is divided into multiple sub-bands, with each sub-band synthesized at a lower sampling rate. These sub-bands are then combined to reconstruct the full-band waveform, thereby accelerating the synthesis process while preserving audio quality.

**CNN-based vocoders.** By leveraging the parallel nature of convolutional operations, CNN-based vocoders [73], [140], [141] can generate high-quality speech more efficiently, making them ideal for real-time applications. A key strength of CNN-based vocoders is their ability to balance synthesis quality and efficiency. However, they often require extensive training data and careful hyperparameter tuning to achieve optimal performance. WaveNet [73] is a probabilistic autoregressive model that generates waveforms sample by sample conditioned on all preceding samples and auxiliary inputs, such as linguistic features and mel-spectrograms. It employs stacks of dilated causal convolutions, enabling long-range dependence modeling in speech signals without relying on recurrent connections. Parallel WaveNet [140] addresses WaveNet's inference speed limitations while maintaining comparable synthesis quality. It introduces a non-autoregressive mechanism based on a teacher-student framework, where the original WaveNet (teacher) distills knowledge into a student model. The student generates samples in parallel, enabling real-time synthesis without waveform quality degradation.

**GAN-based vocoders.** GANs have been widely adopted in vocoders for high-quality speech generation [116], [142]–[145], leveraging adversarial losses to improve realism. GAN-based vocoders typically consist of a generator that produces waveforms conditioned on acoustic features, such as mel-spectrograms, and a discriminator that distinguishes between real and synthesized waveforms. Models like Parallel WaveGAN [144] and HiFi-GAN [116] have demonstrated the effectiveness of GANs in vocoding by introducing tailored loss functions, such as multi-scale and multi-resolution spectrogram losses, to ensure naturalness in both time and frequency domains. These models can efficiently handle the complex, non-linear relationships inherent in speech signals, resulting in high-quality synthesis. A key advantage of GAN-based vocoders is their parallel inference capability, enabling real-time synthesis with lower computational costs compared to autoregressive models. However, training GANs can be

challenging due to instability and mode collapse. Despite these challenges, GAN-based vocoders continue to advance the state-of-the-art in neural vocoding, offering a compelling combination of speed and audio quality.

**Diffusion-based vocoders.** Inspired by diffusion probabilistic models [146] that have shown success in visual generation tasks, diffusion-based vocoders [113], [147]–[150] present a novel approach to natural-sounding speech synthesis. The core mechanism of diffusion-based vocoders involves two stages: a forward process and a reverse process. In the forward process, clean speech waveforms are progressively corrupted by adding noise in a controlled manner, creating a sequence of intermediate noisy representations. During training, the model learns to reverse this process, progressively denoising the corrupted signal to reconstruct the original waveform. Diffusion-based vocoders, such as WaveGrad [149] and DiffWave [148], have demonstrated remarkable performance in generating high-fidelity waveforms while maintaining temporal coherence and natural prosody. They offer advantages over previous vocoders, including robustness to over-smoothing [151] and the ability to model complex data distributions. However, their iterative sampling process can be computationally intensive, posing challenges for real-time applications.

**Other vocoders.** There are also many other types of vocoders such as flow-based [152]–[156] and VAE-based vocoders [157]–[159]. These methods provide unique strengths for speech synthesis such as efficiency and greater flexibility in modeling complex speech variations. Readers can refer to the survey paper from Tan et al. [42] for more details.

The choice of vocoder depends on various factors. While high-quality models like GAN-based and diffusion-based vocoders excel in naturalness, they may not be suitable for real-time scenarios. On the other hand, models like Parallel WaveNet [140] balance quality and efficiency for practical use cases. The best choice will ultimately depend on the specific use case, available resources, and the importance of factors such as model size, training data, and inference speed.

### D. Fully End-to-end TTS models

Fully end-to-end TTS methods [76], [159]–[162] directly generate speech waveforms from textual input, simplifying the "acoustic model → vocoder" pipeline and achieving efficient speech generation. Char2Wav [160] is an early neural text-to-speech (TTS) system that directly synthesizes speech waveforms from character-level text input. It integrates two components and jointly trains them: a recurrent sequence-to-sequence model with attention, which predicts acoustic features (e.g., mel-spectrograms) from text, and a SampleRNN-based neural vocoder [136] that generates waveforms from these features. Similarly, FastSpeech 2s [76] directly synthesizes speech waveforms from texts by extending FastSpeech 2 [76] with a waveform decoder, achieving high-quality and low-latency synthesis. VITS [159] is another fully end-to-end TTS framework. It integrates a variational autoencoder (VAE) with normalizing flows [163] and adversarial training, enabling the model to learn latent representations that capture the intricate variations in speech, such as prosody and style. VITS combines non-autoregressive synthesis with stochastic latent variable modeling, achieving real-time waveform generation without compromising naturalness. There are more end-to-end TTS models such as Tacotron [74], ClariNet [161], and EATS [162], refer to another survey [42] for more details. End-to-end controllable methods that emerged in recent years will be discussed in Section IV.

### E. Acoustic Feature Representations

In TTS, the choice of acoustic feature representations impacts the model's flexibility, quality, expressiveness, and controllability. This subsection investigates continuous representations and discrete tokens as shown in Fig.2, along with their pros and cons for TTS applications.

**Continuous representations.** Continuous representations (e.g., mel-spectrograms) of intermediate acoustic features use a continuous feature space to represent speech signals. These representations often involve acoustic features that capture frequency, pitch, and other characteristics without discretizing the signal. The advantages of continuous features are: 1) Continuous representations retain fine-grained detail, enabling more expressive and natural-sounding speech synthesis. 2) Since continuous features inherently capture variations in tone, pitch, and emphasis, they are well-suited for prosody control and emotional TTS. 3) Continuous representations are more robust to information loss and can avoid quantization artifacts, allowing smoother, less distorted audio. GAN-based [116], [144], [145] and diffusion-based methods [147], [148] often utilize continuous feature representations, i.e., mel-spectrograms. However, continuous representations are typically more computationally demanding and require larger models and memory, especially in high-resolution audio synthesis.

**Discrete tokens.** In discrete token-based TTS, the intermediate acoustic features (e.g., quantized units or phoneme-like tokens) are discrete values, similar to words or phonemes in languages. These are often produced using quantization techniques or learned embeddings, such as HuBERT [166] and SoundStream [168]. The advantages of discrete tokens are: 1) Discrete tokens can encode phonemes or sub-word units, making them concise and less computationally demanding to handle. 2) Discrete tokens often allow TTS systems to require fewer samples to learn and generalize, as the representations are compact and simplified. 3) Using discrete tokens simplifies cross-modal TTS applications like voice cloning or translation-based TTS, as they map well to text-like representations such as LLM tokens. LLM-based [78], [103], [105], [106] and zero-shot TTS methods [17], [78], [87] often adopt discrete tokens as their acoustic features. However, discrete representation learning may result in information loss or lack the nuanced details that can be captured in continuous representations.

Table IV and III summarize the types of acoustic features of representative methods. Table II summarizes popular open-source speech quantization methods.

## III. UNCONTROLLABLE TTS

The development of Uncontrollable Text-To-Speech (UC-TTS) systems represents a significant shift from traditional,

TABLE II
POPULAR OPEN-SOURCE SPEECH QUANTIZATION METHODS.

| Method | Modeling | Code | Year |
|---|---|---|---|
| VQ-Wav2Vec [164] | SSCP | https://github.com/facebookresearch/fairseq/tree/main/examples/wav2vec#vq-wav2vec | 2019 |
| Wav2Vec 2.0 [165] | SSCP | https://github.com/facebookresearch/fairseq/tree/main/examples/wav2vec | 2019 |
| HuBERT [166] | SSCP | https://github.com/facebookresearch/fairseq/tree/main/examples/hubert | 2021 |
| W2v-BERT 2.0 [167] | SSCP | https://huggingface.co/facebook/w2v-bert-2.0 | 2023 |
| SoundStream [168] | RVQGAN | https://github.com/wesbz/SoundStream | 2021 |
| Encodec [169] | RVQGAN | https://github.com/facebookresearch/encodec | 2022 |
| HiFi-Codec [170] | RVQGAN | https://github.com/yangdongchao/AcademiCodec | 2023 |
| SpeechTokenizer [171] | RVQGAN | https://github.com/ZhangXInFD/SpeechTokenizer | 2023 |
| Descript Audio Codec [172] | RVQGAN | https://github.com/descriptinc/descript-audio-codec | 2023 |

SSCP: Self-supervised context prediction, RVQ: Residual vector quantization [168].

linguistics-based synthesis to modern, data-driven deep learning techniques. This shift highlights the integration of both local and global information to produce speech with human-like quality and naturalness. This survey explores UC-TTS evolution, emphasizing the role of local and global information in enhancing speech fidelity and expressiveness.

In the context of UC-TTS, "uncontrollable" refers to the absence of explicit control mechanisms for speech features such as emotion, timbre, and speaker style. Despite this lack of explicit control, the goal is to achieve natural, fluid speech while minimizing issues like mispronunciations and omissions.

### A. Early Approaches: Statistical Models

Early Text-To-Speech (TTS) systems relied on statistical models such as Hidden Markov Models (HMMs) [64], [65] and early neural network-based parametric methods [111], [112]. These models operated at the frame level, using acoustic models and vocoders for text-to-speech conversion. Notable contributions from Tokuda et al. [173] employed HMMs for statistical parametric synthesis, focusing on local features like phonemes, accents, and prosody to improve speech naturalness.

While robust, these statistical methods were limited by their reliance on pre-segmented data, leading to oversimplified assumptions about speech dynamics. Local linguistic features were well-modeled, but the global phonetic context was often overlooked, resulting in speech that sounded monotone and lacked emotional depth, as noted by Zen et al. [174].

### B. Sequence-to-Sequence Models

The emergence of sequence-to-sequence models represents a significant breakthrough by removing the need for explicit linguistic features, thereby enabling the capture of the nuances and idiosyncrasies of human speech. Models such as Tacotron [74] and Tacotron 2 [175] utilize recurrent neural networks (RNNs) with attention mechanisms to effectively model the complex, nonlinear nature of speech sequences. These innovations allow for precise tuning of speech parameters, enhancing prosody and rhythm by modeling entire utterances rather than isolated phonetic units.

Building on these advancements, Deep Voice 3 [176] introduces a fully convolutional sequence-to-sequence architecture that significantly accelerates training speed compared to RNN-based models. This approach achieves training times an order of magnitude faster, enabling scalability to handle large datasets. Additionally, the use of a position-augmented attention mechanism in Deep Voice 3 enhances the naturalness of synthesized speech, achieving competitive mean opinion scores, especially when paired with advanced neural vocoders like WaveNet. This development not only improves training efficiency but also enhances the scalability and naturalness of text-to-speech systems.

### C. Transformer-Based Models

Transformer-based architectures advanced the field by enabling computational parallelization and effectively capturing long-range dependencies. Models like Transformer TTS overcame RNN challenges, such as gradient vanishing, by using efficient training paradigms [177]. Self-attention mechanisms allowed simultaneous modeling of local phonetic details and global prosodic contexts, resulting in more sophisticated and human-like speech synthesis.

Although transformers improved contextual information incorporation, challenges remained in preserving local phonetic precision. To address these, techniques such as relative position encodings and localized attention were integrated [124].

### D. Advanced Architectures: Integrating Flow and Diffusion Models

Recent advancements have shifted towards integrating global information within end-to-end architectures to enhance speech naturalness and coherence. Flow-based models like Glow-TTS [133] and Flow-TTS [132] exemplify this by employing invertible transformations that maintain the balance between local precision and global coherence. These architectures enable the synthesis of high-fidelity speech by modeling complex dependencies across the entire utterance, thus improving the overall fluidity and naturalness of the generated speech.

Moreover, the introduction of diffusion models in TTS, such as WaveGrad 2 [178], highlights the shift towards models that can iteratively refine speech output. These models use score matching and diffusion processes to generate speech directly from phoneme sequences, effectively capturing both local nuances and overarching global patterns. The iterative nature of these models allows for adjustments that enhance the

quality of the synthesized audio, accommodating variations in speech without explicit control over specific attributes.

The integration of adversarial training and variational autoencoders (VAEs) further exemplifies the evolution towards incorporating global information. Systems like VITS [159] leverage these techniques to enhance expressiveness and naturalness by learning complex mappings between text and speech. This approach allows the model to manage variations in prosody and rhythm inherently derived from the textual input, aligning with the objectives of UC-TTS to produce diverse and natural speech outputs.

The evolution from HMMs to advanced architectures in UC-TTS exemplifies progress toward synthesizing speech that is both expressive and precise. The interplay of local and global information is crucial for enhancing speech quality and customizability. Future UC-TTS research aims to produce high-fidelity, customizable speech by harmonizing deep contextual insights with precise local adjustments, meeting diverse user needs and communication contexts.

## IV. CONTROLLABLE TTS

In this section, we first review recent TTS work from the perspective of model architecture, followed by a detailed discussion of control modes in controllable TTS.

### A. Model Architectures

Current model architectures can be broadly classified into two main categories: the first is the non-autoregressive (NAR) generative models, which are based on HMMs, neural networks, VAEs, diffusion models, flow matching, and other NAR techniques. The second category relies on autoregressive (AR) codec language models, which typically quantize speech into discrete tokens and use decoder-only models to autoregressively generate these tokens. We summarize the NAR-based and AR-based controllable TTS methods in Table III and Table IV, respectively.

*1) Fully Non-Autoregressive (NAR) Architectures:*

**HMM-based Approaches.** In the realm of Controllable Text-To-Speech (CTTS), advancements in Hidden Markov Model (HMM) architectures have significantly enhanced the manipulation of speech elements such as emotion and prosody. Yamagishi et al. [70] pioneered this field by introducing style-dependent and style-mixed modeling, which allowed precise emulation of human-like emotional nuances and versatile synthesis across various styles by incorporating style as a contextual variable. Building on this foundation, Qin et al. [179] developed the "average emotion model," which utilized MLLR-based adaptation to modulate emotions like happiness and sadness even with limited data, thus advancing the emotional intelligence of synthetic speech systems.

Furthering expressive variability, Nose et al. [119] integrated subjective style intensities and a multiple-regression global variance model into HMM frameworks, addressing over-smoothing and enabling nuanced emotional expressions. Lorenzo-Trueba et al. [72] expanded on these capabilities with CSMAPLR adaptation, introducing "emotion transplantation" to transfer emotional states between speakers while preserving voice distinctiveness, enhancing personalized human-computer interaction. These innovations in HMM architectures have broadened the expressiveness and individuality in synthetic speech, augmenting technological interfaces and paving the way for future developments in adaptive, lifelike speech solutions.

**Transformer-based Approaches.** Advancements in Controllable Text-to-Speech (TTS) technology highlight the integration of deep learning with audio processing, driven by Transformer-based architectures. Ren et al. [15] introduced FastSpeech, a feed-forward non-autoregressive Transformer model that significantly enhances TTS efficiency by reducing inference time and improving the stability issues found in autoregressive models like Tacotron 2. This model provides precise control over prosodic features through duration prediction, effectively tackling the one-to-many mapping challenge. FastSpeech 2 [180] builds on this by integrating pitch and energy control, eliminating the need for the complex teacher-student distillation process, thus enhancing training efficiency and improving voice quality. Parallel Tacotron [84] further advances TTS by employing a variational autoencoder-based residual encoder, capturing intricate prosodic nuances. This approach, combined with iterative spectrogram loss, significantly enhances the naturalness and quality of synthesized speech. Additionally, FastPitch [77] incorporates direct pitch prediction into its architecture, enabling fully parallelized synthesis and precise pitch manipulation. This capability enhances expressiveness and retains the efficiency benefits established by FastSpeech. These innovations significantly contribute to the development of more interactive and natural AI-driven communication systems, underscoring the potential of integrating AI with human-centric disciplines to craft a future where technology and humanity coexist harmoniously.

**VAE-based Approaches.** Recent advancements in Controllable Text-To-Speech (TTS) systems are largely driven by the integration of Variational Autoencoder (VAE) architectures, which enhance the flexibility and precision of speech modulation. Zhang et al. [18] pioneered the use of VAEs in end-to-end speech synthesis, creating disentangled latent representations that allow effective style control and transfer, especially in prosody and emotion management, outperforming the Global Style Token model in style transfer tasks. Building on this, Hsu et al. [130] developed a hierarchical generative model with a conditional VAE framework and a Gaussian mixture model, enabling precise control over complex speech attributes such as environment and style, thus improving expressive speech synthesis through refined noise and speaker characteristic management. Liu et al. [181] further advanced the field with the CLONE model, a single-stage TTS system that resolves the one-to-many mapping issue and enhances high-frequency information reconstruction. By employing a conditional VAE with normalizing flows and a dual path adversarial training mechanism with multi-band discriminators, CLONE achieves nuanced control over prosody and energy, demonstrating superior performance in both speech quality and prosody control compared to state-of-the-art models. These collective innovations highlight the adaptability of VAEs in managing complex speech generation tasks, marking significant progress toward

more dynamic and versatile TTS technologies, with ongoing research promising even greater advancements.

**Diffusion-based Approaches.** The core concept of diffusion-based models is to generate target data by progressively removing noise. During the forward diffusion phase, noise is incrementally added to the original data to form a noise distribution. In the generation phase, a reverse denoising process is employed to gradually recover high-quality speech from the noise. Grad-StyleSpeech [182] introduces a hierarchical transformer encoder to create a representative noise prior distribution for speaker-adaptive settings using score-based diffusion models. NaturalSpeech 2 [86] uses a neural audio codec with residual vector quantizers to obtain quantized latent vectors, which are then generated using a diffusion model conditioned on text input. NaturalSpeech 3 [87] decomposes speech into distinct subspaces that represent different attributes and generates each subspace independently. DEX-TTS [183] improves DiT-based diffusion networks by applying overlapping patchify and convolution-frequency patch embedding strategies. E3 TTS [184] models the temporal structure of the waveform through the diffusion process, eliminating the need for any intermediate representations such as spectrogram features or alignment information.

Applying diffusion models to TTS requires a complex pipeline due to the need for precise temporal alignment between text and speech and the high fidelity required for audio data. This includes domain-specific modeling, such as phoneme and duration [86]. To address the issue of reduced naturalness caused by the addition of duration models, DiTTo-TTS [185] leverages the off-the-shelf pre-trained text and speech encoders without relying on speech domain-specific modeling by incorporating cross-attention mechanisms with the prediction of the total length of speech representations. Similarly, SimpleSpeech [186] proposes a speech codec model (SQ-Codec) based on scalar quantization and uses the sentence duration to control the generated speech length.

**Flow-based Approaches.** Flow-based methods leverage invertible flow transformations to learn mappings from target speech features to simple distributions, typically standard Gaussian distributions. Due to their invertibility, this mechanism can directly sample from the simple distribution and generate high-fidelity speech in the reverse direction. Audiobox [187] and P-flow [188] employ non-autoregressive flow-matching models for efficient and stable speech synthesis. VoiceBox [189] also employs a flow-matching to generate speech, effectively casting the TTS task into a speech infilling task. SpeechFlow [190] is trained on 60k hours of untranscribed speech with flow matching and mask conditions and can be fine-tuned with task-specific data to match or surpass existing expert models. This highlights the potential of generative models as foundation models for speech applications. HierSpeech++ [191] proposes a hierarchical variational inference method. FlashSpeech [192] is built on a latent consistency model and applies a novel adversarial consistency training approach that can train from scratch without the need for a pre-trained diffusion model as the teacher, achieving speech generation in one or two steps.

Recently, E2 TTS [193] converts text input into a character

sequence with filler tokens and trains a mel spectrogram generator based on audio infilling task, achieving human-level naturalness. Inspired by E2 TTS, F5-TTS [194] refines the text representation with ConvNext v2 [195], facilitating easier alignment with speech. E1 TTS [196] further distills a diffusion-based TTS model into a one-step generator with distribution matching distillation [197], [198], reducing the number of network evaluations in sampling from diffusion models. SimpleSpeech 2 [199] introduces a flow-based scalar transformer diffusion model. The work also provides a theoretical analysis, showing that the inclusion of a small number of noisy labels in a large-scale dataset is equivalent to introducing classifier-free guidance during model optimization.

**Other NAR Approaches.** Other works leverage GAN-based or Masked Generative-based methods for TTS generation. StyleTTS 2 [89] employs large pre-trained speech language models (SLMs) such as Wav2Vec 2.0 [165], Hu-BERT [166], and WavLM [200] as discriminators, in combination with a novel differentiable duration modeling approach. This setup uses SLM representations to enhance the naturalness of the synthesized speech. MaskGCT [78] proposes masked generative transformers without requiring text-speech alignment supervision and phone-level duration. The model employs a two-stage system, both trained using a mask-and-predict learning paradigm.

*2) Autoregressive (AR) Architectures:*

**RNN-based Approaches.** Controllable Text-To-Speech (TTS) technology has seen significant advancements through innovations in neural network architectures, facilitating speech that is both natural-sounding and adaptable in terms of emotion, prosody, and pitch. A key breakthrough was the introduction of Tacotron [74], a sequence-to-sequence architecture that effectively integrates prosodic variations, laying the groundwork for precise control over speech attributes. Tacotron 2 [175] further enhanced this capability by better managing prosodic variability, though it averaged these variations, indicating a need for more sophisticated control methods. To address these constraints, Wang et al. [19] introduced Global Style Tokens (GSTs), using an unsupervised approach to encapsulate diverse speech styles into fixed tokens, thus enabling versatile style transfer within the Tacotron framework. Skerry-Ryan et al. [201] further advanced this by incorporating prosodic embeddings, providing detailed control over timing and intonation, and significantly improving the replication of emotional expressions in synthetic speech.

Building on these innovations, emotion-controllable models developed by Li et al. [21] focus on calibrating emotional nuances using emotion embedding networks and style loss alignment, allowing detailed modulation of emotional strength. Hierarchical models like MsEmoTTS [83] refine this approach by segmenting synthesis into global, utterance-level, and local emotional strengths, offering enhanced emotional expressiveness and intuitive control. These advancements have expanded the scope to produce nuanced TTS outputs, enabling precise control over emotion, prosody, and pitch, with applications ranging from virtual assistants to interactive narratives. As researchers continue to explore the potential of neural networks in TTS, the technology promises even richer, more engaging

digital experiences, moving towards speech synthesis that is indistinguishable from natural human interaction.

**LLM-based Approaches.** Inspired by the success of large language models (LLMs) in natural language processing (NLP), recent studies have explored leveraging in-context learning for zero-shot TTS generation.

VALL-E [85] is a pioneering work in this area, formulating TTS as a conditional language modeling problem. It utilizes EnCodec [202] to discretize waveforms into tokens as intermediate representations and employs a two-stage modeling pipeline: an autoregressive model first generates coarse audio tokens, followed by a non-autoregressive model that iteratively predicts additional codebook codes for refinement. This hierarchical modeling of semantic and acoustic tokens has set the foundation for many subsequent LLM-based TTS approaches [203]–[206].

Building on VALL-E, various improvements have been proposed. VALL-E X [207] extends VALL-E to multilingual scenarios, supporting zero-shot cross-lingual speech synthesis and speech-to-speech translation. ELLA-V [208] introduces a sequence order rearrangement step, enhancing local alignment between phoneme and acoustic modalities. RALL-E [209] incorporates prosody tokens as chain-of-thought prompting [210] to stabilize the generation of speech tokens. VALL-E R [211] improves phoneme-to-acoustic alignment and adopts codec-merging to boost decoding efficiency and reduce computational overhead. VALL-E 2 [212] introduces repetition aware sampling and grouped code modeling for greater stability and faster inference. HALL-E [213] adopts a hierarchical post-training framework, effectively managing the trade-off between reducing frame rate and producing high-quality speech.

Beyond the foundational improvements introduced by VALL-E and its immediate extensions, further advancements have focused on enhancing speech alignment, quality, and robustness. SpearTTS [203] and Make-a-voice [204] use semantic tokens to bridge the gap between text and acoustic features. FireRedTTS [214] further optimizes the tokenizer architecture to enhance speech quality. CoFi-Speech [215] generates speech in a coarse-to-fine manner via a multi-scale speech coding and generation approach, producing natural and intelligible speech. CosyVoice [17] employs supervised semantic tokens to enhance content consistency and speaker similarity in zero-shot voice cloning. Similarly, BASE TTS [216] introduces discrete speech representations based on the WavLM [200] self-supervised model, focusing on phonemic and prosodic information. SeedTTS [217] also proposes a self-distillation method for speech decomposition and a reinforcement learning approach to enhance the robustness, speaker similarity and controllability of generated speech. Based on this framework, Bailing-TTS [218] enhances the alignment of text and speech tokens using a continual semi-supervised learning strategy, enabling high-quality synthesis of Chinese dialect speech.

Although models using discrete tokens as intermediate representations have achieved notable success in zero-shot TTS, they still face fidelity issues compared to the continuous representation like Mel spectrograms [219], [220]. MELLE [219] optimizes the training objectives and sampling strategy, marking the first exploration of using continuous-valued tokens instead of discrete-valued tokens within the paradigm of autoregressive speech synthesis models. Similar to MELLE, ARDiT [220] encodes audio as vector sequence in continuous space and autoregressively generates these sequences by a decoder-only transformer.

## B. Control Strategies

The control strategies in existing controllable TTS can be broadly classified into four categories: style tagging using discrete labels, speech reference prompt for customizing a new speaker's voice with just a few seconds of voice input, controlling speech style using natural language descriptions, and the instruction-guided mode. We illustrate taxonomies of controllable TTS from the perspective of control strategies in Fig. 3.

*1) Style Tagging:* This paradigm typically employs target control attributes, primarily emotion-related controls, as categorical label inputs to enable controllable speech synthesis. StyleTagging-TTS [223] utilizes a short phrase or word to represent the style of an utterance and learns the relationship between linguistic embedding and style embedding space by a pre-trained language model.

However, these methods are limited in expressive diversity, as they can only model a small set of pre-defined styles.

*2) Reference Speech Prompt:* This paradigm aims to customize a new speaker's voice with just a few seconds of voice prompt. The architecture can be abstracted into two main components: a speaker encoder that processes the reference speech and outputs a speaker embedding, and a conditional TTS decoder that takes both text and speaker embedding as input to generate speech that matches the style of the reference prompt. MetaStyleSpeech [225] and StyleTTS [88] use adaptive normalization as a style conditioning method, enabling robust zero-shot performance. GenerSpeech [90] introduces a multi-level style adapter to improve zero-shot style transfer for out-of-domain custom voices. SC VALL-E [205] facilitates control over synthesized speech's emotions, speaking styles, and various acoustic features by incorporating style tokens and scale factors. ArtSpeech [236] revisits the sound production system by integrating articulatory representations into the TTS framework, improving the physical interpretability of articulation movements.

To enhance the learning of contextual information and address the challenge of limited voice data from the target speaker, CCSP [237] proposes a contrastive context-speech pretraining framework that learns cross-modal representations combining both contextual text and speech expressions. DEX-TTS [183] separates styles into time-invariant and time-variant components, enabling the extraction of diverse styles from expressive reference speech. StyleTTS-ZS [239] leverages distilled time-varying style diffusion to capture diverse speaker identities and prosodies.

Some works also decouple timbre and style information from the reference speech, allowing more flexible control over the speaking style [87], [106], [230]. MegaTTS 2 [230] introduces an acoustic autoencoder that separately encodes prosody

TABLE III
A SUMMARY OF EXISTING NON-AUTOREGRESSIVE CONTROLLABLE NEURAL-BASED METHODS.

| Method | Zero-shot TTS | Controlability | | | | | | | | Model Architectures | | Acoustic Feature | Release Time |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Pit. | Ene. | Spe. | Pro. | Tim. | Emo. | Env. | Des. | Acoustic Model | Vocoder | | |
| FastSpeech [15] | | | | ✓ | ✓ | | | | | Transformer | WaveGlow | MelS | 2019.05 |
| DWAPI [221] | | ✓ | | | ✓ | ✓ | | | | DNN | Straight | MelS + F0 + Intensity | 2020.04 |
| FastSpeech 2 [180] | | ✓ | ✓ | ✓ | ✓ | | | | | Transformer | Parallel WaveGAN | MelS | 2020.06 |
| FastPitch [77] | | ✓ | | | ✓ | | | | | Transformer | WaveGlow | MelS | 2020.06 |
| Parallel Tacotron [222] | | | | | ✓ | | | | | Transformer + CNN | WaveRNN | MelS | 2020.10 |
| StyleTagging-TTS [223] | ✓ | | | | ✓ | ✓ | | | | Transformer + CNN | HiFi-GAN | MelS | 2021.04 |
| SC-GlowTTS [224] | ✓ | | | | ✓ | | | | | Transformer + Conv | HiFi-GAN | MelS | 2021.06 |
| Meta-StyleSpeech [225] | ✓ | | | | ✓ | | | | | Transformer | MelGAN | MelS | 2021.06 |
| DelightfulTTS [226] | | ✓ | | ✓ | ✓ | | | | | Transformer + CNN | HiFiNet | MelS | 2021.11 |
| YourTTS [82] | ✓ | | | | ✓ | | | | | Transformer | HiFi-GAN | LinS | 2021.12 |
| DiffGAN-TTS [227] | | ✓ | | ✓ | ✓ | | | | | Diffusion + GAN | HiFi-GAN | MelS | 2022.01 |
| StyleTTS [88] | ✓ | | | | ✓ | | | | | CNN + RNN + GAN | HiFi-GAN | MelS | 2022.05 |
| GenerSpeech [90] | ✓ | | | | ✓ | | | | | Transformer + Flow-based | HiFi-GAN | MelS | 2022.05 |
| NaturalSpeech 2 [86] | ✓ | | | | ✓ | | | | | Diffusion | Codec Decoder | Token | 2022.05 |
| Cauliflow [228] | | | | ✓ | ✓ | | | | | BERT + Flow | UP WaveNet | MelS | 2022.06 |
| CLONE [181] | | ✓ | | ✓ | ✓ | | | | | Transformer + CNN | WaveNet | MelS + LinS | 2022.07 |
| PromptTTS [101] | | ✓ | ✓ | ✓ | ✓ | | ✓ | | ✓ | Transformer | HiFi-GAN | MelS | 2022.11 |
| Grad-StyleSpeech [182] | ✓ | | | | ✓ | | | | | Score-based Diffusion | HiFi-GAN | MelS | 2022.11 |
| PromptStyle [229] | ✓ | ✓ | | | ✓ | ✓ | ✓ | | ✓ | VITS | HiFi-GAN | MelS | 2023.05 |
| StyleTTS 2 [89] | ✓ | | | | ✓ | ✓ | ✓ | | | Diffusion + GAN | HifiGAN / iSTFTNet | MelS | 2023.06 |
| VoiceBox [189] | ✓ | | | | ✓ | | | | | Flow Matching Diffusion | HiFi-GAN | MelS | 2023.06 |
| MegaTTS 2 [230] | ✓ | | | | ✓ | ✓ | | | | Diffusion + GAN | HiFi-GAN | MelS | 2023.07 |
| PromptTTS 2 [102] | | ✓ | ✓ | ✓ | ✓ | | | | ✓ | Diffusion | Codec Decoder | Token | 2023.09 |
| VoiceLDM [231] | | ✓ | | | ✓ | ✓ | ✓ | ✓ | ✓ | Diffusion | HiFi-GAN | MelS | 2023.09 |
| DuIAN-E [232] | | ✓ | | ✓ | ✓ | | | | | CNN + RNN | HiFi-GAN | MelS | 2023.09 |
| PromptTTS++ [233] | | ✓ | | ✓ | ✓ | ✓ | ✓ | | ✓ | Transformer + Diffusion | BigVGAN | MelS | 2023.09 |
| SpeechFlow [190] | ✓ | | | | ✓ | | | | | Flow Matching Diffusion | HiFi-GAN | MelS | 2023.10 |
| P-Flow [188] | ✓ | | | | ✓ | | | | | Flow Matching | HiFi-GAN | MelS | 2023.10 |
| E3 TTS [184] | ✓ | | | | ✓ | | | | | Diffusion | / | Waveform→Unet | 2023.11 |
| HierSpeech++ [191] | ✓ | | | | ✓ | | | | | Hierarchical Conditional VAE | BigVGAN | MelS | 2023.11 |
| Audiobox [187] | ✓ | ✓ | | ✓ | ✓ | | | ✓ | ✓ | Flow Matching | EnCodec | MelS | 2023.12 |
| FlashSpeech [192] | ✓ | | | | ✓ | | | | | Latent Consistency Model | EnCodec | Token | 2024.04 |
| NaturalSpeech 3 [87] | ✓ | | | ✓ | ✓ | ✓ | | | | Diffusion | EnCodec | Token | 2024.04 |
| InstructTTS [105] | | ✓ | | ✓ | ✓ | ✓ | ✓ | | ✓ | Transformer + Diffusion | HiFi-GAN | Token | 2024.05 |
| ControlSpeech [106] | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | ✓ | Transformer + Diffusion | FACodec Decoder | Token | 2024.06 |
| AST-LDM [234] | | | | | ✓ | | | ✓ | ✓ | Diffusion | HiFi-GAN | MelS | 2024.06 |
| SimpleSpeech [186] | ✓ | | | | ✓ | | | | | Transformer Diffusion | SQ Decoder | Token | 2024.06 |
| DiTTo-TTS [185] | ✓ | | | ✓ | ✓ | | | | | DiT | BigVGAN | Token | 2024.06 |
| E2 TTS [193] | ✓ | | | | ✓ | | | | | Flow Matching Transformer | BigVGAN | MelS | 2024.06 |
| MobileSpeech [235] | ✓ | | | | ✓ | | | | | ConFormer Decoder | Vocos | Token | 2024.06 |
| DEX-TTS [183] | ✓ | | | | ✓ | | | | | Diffusion | HiFi-GAN | MelS | 2024.06 |
| ArtSpeech [236] | ✓ | | | | ✓ | | | | | RNN + CNN | HiFI-GAN | MelS+Energy+F0+TV | 2024.07 |
| CCSP [237] | ✓ | | | | ✓ | | | | | Diffusion | Codec Decoder | Token | 2024.07 |
| SimpleSpeech 2 [199] | ✓ | | | ✓ | ✓ | | | | | Flow-based Transformer Diffusion | SQ Decoder | Token | 2024.08 |
| E1 TTS [196] | ✓ | | | | ✓ | | | | | DiT | BigVGAN | Continuous Token | 2024.09 |
| VoiceGuider [238] | ✓ | | | | ✓ | | | | | Diffusion | BigVGAN | MelS | 2024.09 |
| StyleTTS-ZS [239] | ✓ | | | | ✓ | | | | | Diffusion + GAN | HifiGAN / iSTFTNet | Token | 2024.09 |
| NansyTTS [240] | ✓ | ✓ | | ✓ | ✓ | ✓ | | | ✓ | Transformer | NANSY++ | MelS | 2024.09 |
| NanoVoice [241] | ✓ | | | | ✓ | | | | | Diffusion | BigVGAN | MelS | 2024.09 |
| MS²KU-VTTS [242] | | | | | | | | ✓ | ✓ | Diffusion | BigVGAN | MelS | 2024.10 |
| MaskGCT [78] | ✓ | | | ✓ | ✓ | | | | | Masked Generative Transformers | DAC + Vocos | Token | 2024.10 |

Abbreviations: Pit(ch), Ene(rgy), Spe(ed), Pro(sody), Tim(bre), Emo(tion), Env(ironment), Des(cription). Timbre involves gender and age. MelS and LinS represent Mel Spectrogram and Linear Spectrogram respectively. TV represents the vocal tract variables proposed in [236].

and timbre into the latent space, enabling the transfer of various speaking styles to the desired timbre. ControlSpeech [106] uses bidirectional attention and mask-based parallel decoding to capture codec representations in a discrete decoupling codec space, allowing independent control of timbre, style, and content in a zero-shot manner.

*3) Natural Language Descriptions:* Recent studies explore controlling speech style using natural language descriptions that include attributes such as pitch, gender, and emotion, making the process more user-friendly and interpretable. In this paradigm, several speech datasets with natural language descriptions [101], [106], [247] and associated prompt generation pipelines [102], [247], [252] have been proposed. Detailed information about these datasets will be discussed in Section V. PromptTTS [101] uses manually annotated text prompts to describe five speech attributes, including gender,

pitch, speaking speed, volume, and emotion. InstructTTS [105] introduces a three-stage training procedure to capture semantic information from natural language style prompts and adds further annotation to the NLSpeech dataset's speech styles. PromptStyle [229] constructs a shared space for stylistic and semantic representations through a two-stage training process. TextrolSpeech [247] proposes an efficient prompt programming methodology and a multi-stage discrete style token-guided control framework, demonstrating strong in-context capabilities. NansyTTS [240] combines a TTS trained on the target language with a description control model trained on another language, which shares the same timbre and style representations to enable cross-lingual controllability.

Considering not all details about voice variability can be described in the text prompt, PromptTTS++ [233] and PromptSpeaker [251] tries to construct text prompts with

TABLE IV
A SUMMARY OF EXISTING AUTOREGRESSIVE CONTROLLABLE NEURAL-BASED METHODS.

| Method | Zero-shot TTS | Controlability | | | | | | | | Model Architectures | | Acoustic Feature | Release Time |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Pit. | Ene. | Spe. | Pro. | Tim. | Emo. | Env. | Des. | Acoustic Model | Vocoder | | |
| Prosody-Tacotron [201] | | ✓ | | | ✓ | | | | | RNN | WaveNet | MelS | 2018.03 |
| GST-Tacotron [243] | | ✓ | | | ✓ | | | | | CNN + RNN | Griffin-Lim | LinS | 2018.03 |
| GMVAE-Tacotron [130] | | ✓ | ✓ | ✓ | | | | ✓ | | CNN + RNN | WaveRNN | MelS | 2018.12 |
| VAE-Tacotron [18] | | ✓ | ✓ | ✓ | | | | | | CNN + RNN | WaveNet | MelS | 2019.02 |
| DurIAN [244] | | ✓ | ✓ | ✓ | | | | | | CNN + RNN | MB-WaveRNN | MelS | 2019.09 |
| Flowtron [245] | | ✓ | ✓ | ✓ | | | | | | CNN + RNN | WaveGlow | MelS | 2020.07 |
| MsEmoTTS [83] | | ✓ | | ✓ | | | ✓ | | | CNN + RNN | WaveRNN | MelS | 2022.01 |
| VALL-E [85] | ✓ | | | | | ✓ | | | | LLM | EnCodec | Token | 2023.01 |
| SpearTTS [203] | ✓ | | | | | ✓ | | | | LLM | SoundStream | Token | 2023.02 |
| VALL-E X [207] | ✓ | | | | | ✓ | | | | LLM | EnCodec | Token | 2023.03 |
| Make-a-voice [204] | ✓ | | | | | ✓ | | | | LLM | BigVGAN | Token | 2023.05 |
| TorToise [246] | | | | | | ✓ | | | | Transformer + DDPM | Univnet | MelS | 2023.05 |
| MegaTTS [91] | ✓ | | | | | ✓ | | | | LLM + GAN | HiFi-GAN | MelS | 2023.06 |
| SC VALL-E [205] | ✓ | ✓ | | ✓ | ✓ | ✓ | ✓ | | | LLM | EnCodec | Token | 2023.07 |
| Salle [247] | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | ✓ | LLM | Codec Decoder | Token | 2023.08 |
| UniAudio [248] | ✓ | ✓ | | ✓ | ✓ | ✓ | | | ✓ | LLM | EnCodec | Token | 2023.10 |
| ELLA-V [208] | ✓ | | | | | ✓ | | | | LLM | EnCodec | Token | 2024.01 |
| BaseTTS [216] | ✓ | | | | | ✓ | | | | LLM | UnivNet | Token | 2024.02 |
| ClaM-TTS [249] | ✓ | | | | | ✓ | | | | LLM | BigVGAN | MelS+Token | 2024.04 |
| RALL-E [209] | ✓ | | | | | ✓ | | | | LLM | SoundStream | Token | 2024.05 |
| ARDiT [220] | ✓ | | | ✓ | | ✓ | | | | Decoder-only Diffusion Transformer | BigVGAN | MelS | 2024.06 |
| VALL-E R [211] | ✓ | | | | | ✓ | | | | LLM | Vocos | Token | 2024.06 |
| VALL-E 2 [212] | ✓ | | | | | ✓ | | | | LLM | Vocos | Token | 2024.06 |
| Seed-TTS [217] | ✓ | | | | | ✓ | ✓ | | | LLM + Diffusion Transformer | / | Token | 2024.06 |
| VoiceCraft [93] | ✓ | | | | | ✓ | | | | LLM | HiFi-GAN | Token | 2024.06 |
| XTTS [250] | ✓ | | | | | ✓ | | | | LLM + GAN | HiFi-GAN | MelS+Token | 2024.06 |
| CosyVoice [17] | ✓ | ✓ | | ✓ | ✓ | ✓ | ✓ | | ✓ | LLM + Conditional Flow Matching | HiFi-GAN | Token | 2024.07 |
| MELLE [219] | ✓ | | | | | ✓ | | | | LLM | HiFi-GAN | MelS | 2024.07 |
| Bailing TTS [218] | ✓ | | | | | ✓ | | | | LLM + Diffusion Transformer | / | Token | 2024.08 |
| VoxInstruct [103] | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | ✓ | LLM | Vocos | Token | 2024.08 |
| Emo-DPO [31] | | | | | | | ✓ | | ✓ | LLM | HiFi-GAN | Token | 2024.09 |
| FireRedTTS [214] | ✓ | | | | ✓ | ✓ | | | | LLM + Conditional Flow Matching | BigVGAN-v2 | Token | 2024.09 |
| CoFi-Speech [215] | ✓ | | | | | ✓ | | | | LLM | BigVGAN | Token | 2024.09 |
| Takin [206] | ✓ | ✓ | | ✓ | ✓ | ✓ | ✓ | | ✓ | LLM | HiFi-Codec | Token | 2024.09 |
| HALL-E [213] | ✓ | | | | | ✓ | | | | LLM | EnCodec | Token | 2024.10 |

Abbreviations: Pit(ch), Ene(rgy), Spe(ed), Pro(sody), Tim(bre), Emo(tion), Env(ironment), Des(cription). Timbre involves gender and age. MelS and LinS represent Mel Spectrogram and Linear Spectrogram respectively.
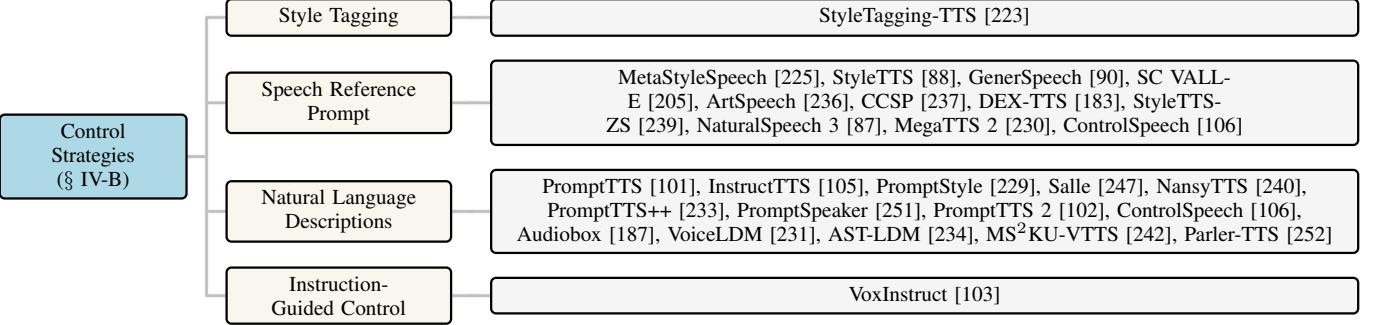


Fig. 3. A taxonomy of controllable TTS from the perspective of control strategies.

more details. PromptTTS 2 [102] designs a variation network to capture voice variability not conveyed by text prompts. ControlSpeech [106] proposes the Style Mixture Semantic Density (SMSD) module, incorporating a noise perturbation mechanism to tackle the many-to-many problem in style control and enhance style diversity.

Other works also focus on improving controllability in additional aspects, such as the surrounding environment. Audiobox [187] introduces both description-based and example-based prompting, integrating speech and sound generation paradigms to independently control transcript, vocal, and other audio styles during speech generation. VoiceLDM [231] and AST-LDM [234] extend AudioLDM [253] to incorporate environmental context in TTS generation by adding a content prompt as a conditional input. Building on VoiceLDM, MS$^2$KU-VTTS [242] further expands the dimensions of environmental perception, enhancing the generation of immersive spatial speech.

*4) Instruction-Guided Control:* The natural language description-based TTS methods discussed above require splitting inputs into content and description prompts, which limits fine-grained control over speech and does not align with other AIGC models. VoxInstruct [103] proposes a new paradigm

TABLE V
A SUMMARY OF OPEN-SOURCE DATASETS FOR CONTROLLABLE TTS.

| Dataset | Hours | #Speakers | Pit. | Ene. | Spe. | Age | Gen. | Emo. | Emp. | Acc. | Top. | Des. | Env. | Dia. | Lang | Release Time |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Taskmaster-1 [254] | / | / | | | | | | | | | | | | ✓ | en | 2019.09 |
| Libri-light [255] | 60,000 | 9,722 | | | | | | | | | ✓ | | | | en | 2019.12 |
| AISHELL-3 [256] | 85 | 218 | | | | ✓ | ✓ | | | ✓ | | | | | zh | 2020.10 |
| ESD [257] | 29 | 10 | | | | | | ✓ | | | | | | | en,zh | 2021.05 |
| GigaSpeech [258] | 10,000 | / | | | | | | | | | ✓ | | | | en | 2021.06 |
| WenetSpeech [259] | 10,000 | / | | | | | | | | | ✓ | | | | zh | 2021.07 |
| PromptSpeech [101] | / | / | ✓ | ✓ | ✓ | | | ✓ | | | | ✓ | | | en | 2022.11 |
| DailyTalk [260] | 20 | 2 | | | | | | ✓ | | | ✓ | | | ✓ | en | 2023.05 |
| TextrolSpeech [247] | 330 | 1,324 | ✓ | ✓ | ✓ | | ✓ | ✓ | | | | ✓ | | | en | 2023.08 |
| VoiceLDM [231] | / | / | ✓ | | | | ✓ | ✓ | | | | ✓ | ✓ | | en | 2023.09 |
| VccmDataset [106] | 330 | 1,324 | ✓ | ✓ | ✓ | | ✓ | ✓ | | | | ✓ | | | en | 2024.06 |
| MSceneSpeech [261] | 13 | 13 | | | | | | | | | ✓ | | | | zh | 2024.07 |
| SpeechCraft [262] | 2,391 | 3,200 | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | ✓ | ✓ | | | en,zh | 2024.08 |

Abbreviations: Pit(ch), Ene(rgy)=volume=loudness, Spe(ed)=duration, Gen(der), Emo(tion), Emp(hasis), Acc(ent), Dia(logue), Env(ironment), Des(cription).

TABLE VI
COMMON OBJECTIVE AND SUBJECTIVE EVALUATION METRICS

| Metric | Type | Eval Target | GT Required |
|---|---|---|---|
| MCD [263] | Objective | Acoustic similarity | ✓ |
| PESQ [264] | Objective | Perceptual quality | ✓ |
| WER [265] | Objective | Intelligibility | ✓ |
| MOS [266] | Subjective | Preference | |
| CMOS [267] | Subjective | Preference | |

GT: Ground truth.

that extends traditional text-to-speech tasks into a general human instruction-to-speech task. Here, human instructions are freely written in natural language, encompassing both the spoken content and descriptive information about the speech. To enable automatic extraction of the synthesized speech content from raw text instructions, VoxInstruct uses speech semantic tokens as an intermediate representation, bridging the gap in current research by allowing the simultaneous use of both text description prompts and speech prompts for speech generation.

## V. DATASETS AND EVALUATION

### A. Datasets

Achieving fully controllable TTS requires large-scale datasets with rich diversity and fine-grained annotations. In this subsection, we categorize speech datasets into four types according to the labels they provide, i.e., attribute tags such as age and gender, description, environment, and dialogues.

**Tag-based speech datasets.** Tag-based datasets [255]–[261] are specialized collections of speech data that include metadata tags, such as pitch, emotion, and energy, age, gender, and so on, to enhance the expressiveness and control of TTS systems. These datasets not only provide audio-text pairs but also offer additional labels that guide the model in generating diverse and context-aware speech outputs.

**Description-based speech datasets.** Description-based TTS datasets [101], [106], [231], [247], [262] contain speech data paired with detailed textual descriptions of the speech

attributes or characteristics, such as emotion, speed, intonation, and style. These datasets enable training TTS models that interpret descriptive prompts to generate expressive and context-aware speech.

**Speech environment datasets.** Speech environment datasets [231], [268] are collections of speech data annotated with environmental labels, such as park, stadium, office, or street, capturing the acoustic characteristics of various real-world settings. These datasets are crucial for training models to handle diverse acoustic scenarios, improving their robustness and adaptability.

**Dialogue speech datasets.** Dialogue speech datasets [254], [260] capture conversational speech between two or more people, focusing on the natural flow, turn-taking, and context of human dialogue. These datasets are crucial for training systems in applications like chatbots, virtual assistants, and interactive storytelling.

We summarize open-source speech datasets for controllable TTS in Table V.

### B. Evaluation

The performance of controllable TTS often requires objective and subjective evaluation. We introduce common evaluation metrics in this subsection.

**Objective evaluation metrics.** Objective metrics offer automated and reproducible evaluations. Mel Cepstral Distortion (MCD) [263] measures the spectral distance between synthesized and reference speech, reflecting how closely the generated audio matches the target in terms of acoustic features. For intelligibility, the Word Error Rate (WER) [265] is used, comparing transcriptions of synthesized speech to the input text via automated speech recognition. Perceptual Evaluation of Speech Quality (PESQ) [264] is another objective metric designed to evaluate speech quality by comparing degraded audio with a clean reference. It is widely used in telecommunications and speech synthesis, PESQ models human auditory perception, producing a score (typically 1–4.5) that reflects intelligibility and distortion under various conditions, including noise or compression.

**Subjective evaluation metrics.** The Mean Opinion Score (MOS) [266] is the most commonly used subjective metric. In MOS evaluations, listeners rate the naturalness of synthesized speech on a scale (e.g., 1 to 5), where higher scores indicate better quality. MOS captures human perception effectively but requires substantial resources for large-scale evaluations. Comparison Mean Opinion Score (CMOS) [267] further evaluates relative quality differences between two TTS audio samples. Participants listen to paired samples and rate their preference on a scale (e.g., -3 to +3, where negative values favor the first sample). CMOS is used to measure subtle improvements in TTS systems, complementing absolute MOS ratings.

Table VI summarizes common evaluation metrics for TTS.

## VI. CHALLENGES AND FUTURE DIRECTIONS

In this section, we elaborate on current challenges for fully controllable TTS and discuss promising future directions.

### A. Challenges

Controllable TTS aims to synthesize speech while allowing precise control over speech characteristics such as pitch, duration, energy, prosody, speaking style, and emotion. While significant progress has been made, achieving truly controllable TTS remains a complex task due to the multifaceted nature of human speech and the technical challenges in modeling and synthesizing it. In this section, we delve into the primary challenges and analyze their underlying reasons.

**Controllable granularity.** A critical challenge in controllable TTS is determining what aspects of speech should be controlled and how to control speech characteristics at a specific granularity. Different applications require varying levels of control granularity. For instance, audiobook narration may need sentence-level control of emotion, while conversational AI like ChatGPT may require word or phoneme-level control over prosody. Moreover, the emotion, prosody, and other characteristics of human speech are often intricately intertwined and can manifest across varying levels of granularity. Additionally, achieving fine-grained control requires high-resolution annotations and sophisticated models capable of handling subtle variations without compromising synthesis quality.

Although some LLM-based TTS methods such as VoxInstruct [103] can control various aspects of speech through attribute descriptions, determining the appropriate level of granularity for control and devising methods to achieve precise control at a *specific granularity* or to enable *multiscale and fine-grained control* remains a significant challenge.

**Feature extraction and representation.** Achieving fully controllable TTS needs good feature disentanglement. Accurately extracting meaningful and disentangled speech features like pitch contours, energy patterns, emotion variation, and prosodic elements from training data is difficult. The reason is that speech features are interdependent and context-sensitive, making it hard to isolate specific attributes for control. For example, altering pitch often affects prosody, emotion, and naturalness to some extent. To tackle this, several methods [269]–[271] utilize pre-trained models for different

speech recognition tasks (e.g., pitch, energy, and duration prediction, gender classification, age estimation, and speaker verification) to supervise feature extraction. For example, NaturalSpeech3 [14] factorizes speech into separate feature subspaces to capture different speech attributes.

However, these methods are limited to coarse or high-level feature disentanglement, leaving a significant gap in *fully disentangled control*. On the other hand, selecting *suitable representations* (e.g., continuous variables like mel-spectrograms or latent embeddings like tokens) for controllable attributes is non-trivial because representations must be both interpretable for humans and expressive enough for TTS models. For example, transformer-based models are good at processing discrete tokens, while GAN and Diffusion-based models excel in modeling continuous representations.

**Scarcity of datasets.** High-quality, diverse, and appropriately annotated datasets are essential for training controllable TTS systems. However, such datasets are scarce and difficult to construct. To achieve controllable TTS, training data must encompass a wide range of styles, emotions, accents, and prosodic variations to enable versatile control because limited diversity in datasets can restrict the model's ability to generalize across unseen styles or emotions. Although there are some TTS datasets, such as LibriTTS [272], Gigaspeech [258], and TextrolSpeech [247], their diversity is still not enough for fully controllable TTS due to the lack of corpus of *diverse content* such as comedies, thrillers, cartoons, etc. Constructing large-scale datasets with rich diversity is also expensive and time-consuming.

Another obstacle is that creating datasets with fine-grained, attribute-specific annotations is labor-intensive and costly. Besides, manual annotation of speech attributes requires expert knowledge and is prone to inconsistencies and errors, particularly for subjective qualities like emotion. Currently, most datasets provide only coarse labels, such as gender, age, or a limited range of emotions. While some datasets, such as SpeechCraft [262], include natural language descriptions of speech attributes, no existing dataset offers *fine-grained variations and annotations* within the speech of the same speakers. Available datasets for controllable TTS are summarized in Table V.

**Generalization ability.** The ability of a TTS system to generalize effectively is crucial for producing natural, high-quality speech across a wide range of conditions, such as unseen speakers, languages, or topics. However, achieving robust generalization remains a significant challenge for modern TTS methods due to various factors. *Zero-shot controllable* TTS [78], [92] aims to synthesize speech for unseen speakers with various speech customization such as emotion using minimal reference audio, which can offer flexibility for personalized voice generation. However, it faces significant challenges, including capturing unique speaker characteristics from limited data, accurately reproducing prosody and style, and disentangling speaker identity from other audio attributes like emotion or noise.

*Multilingual generalization* [250], [273] in TTS refers to the ability to synthesize natural and intelligible speech across multiple languages, including those not seen during training.

This capability is essential for applications like cross-lingual communication, multilingual virtual assistants, and speech synthesis for low-resource languages [274]. Multilingual generalization still faces many challenges such as linguistic diversity and mismatch and the scarcity of data. Cross-lingual speaker generalization is another hurdle, as preserving speaker identity across languages can lead to artifacts.

*Domain adaptation* [275] in TTS refers to tailoring a pretrained TTS model to generate speech for a specific domain or context, such as medical terminology and conversational speech. One challenge is that many specialized domains lack sufficient high-quality annotated data for fine-tuning. Besides, adapting prosody, intonation, and speaking style to match domain-specific requirements such as comic dialogue is complex. Failing to capture domain-specific nuances can make speech sound unnatural or inconsistent with the target context.

**Efficiency.** Efficiency in controllable TTS systems is a critical requirement for practical applications, as these models aim to offer fine-grained control over various speech attributes such as prosody, emotion, style, and speaker identity. However, achieving such control often comes at the cost of increased computational complexity, larger model sizes, and longer inference times, creating significant challenges.

High latency is a major issue, as existing controllable TTS models [78], [101]–[103] often necessitate autoregressive processes to synthesize speech. The inference time of these models ranges from several to tens of seconds. This can be especially problematic for real-time applications like live broadcasting or interactive systems. Additionally, the challenge of balancing granularity and efficiency arises, as finer controls demand higher-resolution data and more precise models, leading to increased resource requirements and *inefficient training and inference*.

Another major obstacle lies in the trade-off between model complexity and performance. State-of-the-art controllable TTS systems often rely on large neural networks such as LLMs with billions of parameters, which provide superior naturalness and expressiveness but demand significant computational resources. Simplifying these architectures can lead to quality degradation, including artifacts, unnatural prosody, or limited expressiveness. Therefore, designing *light-weight* controllable TTS models is significantly tricky.

### B. Future Directions

In this survey, we conduct a comprehensive investigation and analysis of existing TTS methods, particularly on controllable TTS technologies. While these methods show great potential in real-world applications, there are still some limitations that need to be addressed. Based on our observations, we outline several promising future directions as follows:

**Fine-grained speech synthesis by natural language description.** Using natural language description to synthesize human speech with fine-grained control over various audio attributes is currently underexplored. Most of the existing works can only control a fixed number of attributes of the synthesized speech. Although a few works show great control of emotion, timbres, pitch, gender, and styles, e.g., VoxInstruct [103] and
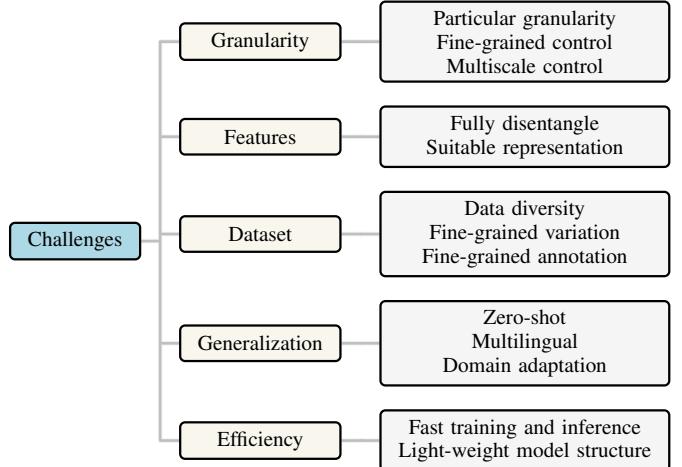


Fig. 4. A summary of current major challenges for fully controllable TTS.

CosyVoice [17], they can frequently synthesize unwanted speech clips. Users need to synthesize multiple times to get satisfactory speech.

**Fine-grained speech editing by natural language description.** Speech or audio editing has been studied for a long time. However, existing methods usually train conditional models and adjust a fixed number of conditional inputs to modify the attributes of synthesized speech, thus lacking fine-grained manipulations [94], [95]. Therefore, how to learn disentangled speech representations for speech attributes while supporting editing by using natural language description is worthy of investigation.

**Expressive multi-modal speech synthesis.** Synthesizing speech from multi-modal data such as texts, images, and videos is an appealing research topic due to its various applications in the industry such as storytelling, filming, and gaming. Although there are several related works on this task [6], [24], [276], [277], few of them can fully extract useful information from multi-modal data. Particularly, synthesizing engaging speech and expressive voiceover for complex visual content sees great opportunities in the future.

**Natural and emotional conversational TTS.** Speech conversational TTS have come out for several decades but remained as cascaded systems for a long time and cannot generate natural and emotional speech. These systems are not context-aware, making the synthesized speech sound robotic. With the advent of LLMs, existing TTS technologies were directly introduced by simply synthesizing speech from the text generated by LLMs [33]. However, context-aware conversational TTS with rich emotion and good naturalness has not been well studied.

**Zero-shot long speech synthesis with emotion consistency.** Zero-shot TTS emerged in recent years to achieve voice cloning and speech style imitation without fine-tuning, making them more practical in real scenarios [17], [78], [194]. However, synthesizing long speech with rich emotion and style variation in a zero-shot setting remains challenging due to the lack of rich speech information in short reference audio clips. Addressing this issue will make a big step towards fully

controllable zero-shot TTS.

**Efficient TTS by natural language description.** Synthesizing speech with natural language description usually involves training large language encoders and bridge nets between the two modalities which can bring about much more computation overhead compared to previous TTS methods. The inference time is also relatively slow, e.g., existing methods usually take tens of seconds to synthesize a short speech audio clip of less than 10 seconds [17], [104]. Therefore, efficient text and speech modeling and interaction is critical for natural language description-based TTS systems.

## VII. DISCUSSION

### A. Impacts of Controllable TTS

**Applications.** Controllable TTS systems allow fine-grained manipulation of speech attributes such as pitch, emotion, speaking rate, and style, enabling a wide range of applications across industries. One major application is in virtual assistants and customer support systems, where controllable TTS ensures tailored and context-aware responses. For instance, a virtual assistant can speak in a calm tone for technical support but switch to an enthusiastic tone when presenting promotional offers, enhancing user experience.

In the entertainment industry, controllable TTS is invaluable for creating dynamic voiceovers, audiobooks, and gaming characters. It enables precise adjustments in tone and delivery, allowing audiobooks to reflect character emotions and gaming characters to exhibit personality traits that align with their roles. Similarly, in education, TTS systems can adapt speaking styles to suit different learners, such as adopting a slow, clear tone for language learning or an engaging, storytelling style for children's content.

Controllable TTS is also transformative in assistive technologies, where it can generate speech that reflects the user's intended emotion or personality. This is particularly impactful for individuals with speech impairments, enabling more expressive and natural communication.

In content localization, controllable TTS systems adjust speech characteristics to suit regional and cultural preferences, ensuring a natural fit for global audiences. Additionally, it finds applications in human-computer interaction research, enabling adaptive dialogue systems that modify speech dynamically based on user mood or environment. By offering this flexibility, controllable TTS systems enhance accessibility, personalization, and engagement across various domains.

**Deepfakes.** A deepfake is a type of synthetic media in which a person in an existing image, video, or audio recording is replaced with someone else's likeness or voice. This technology uses deep learning, particularly GANs [278], to create highly realistic, but fabricated, content. While deepfakes are most commonly associated with video manipulation, such as face swapping [279], they can also apply to audio, enabling the creation of synthetic speech that mimics a specific person's voice, which is well known as voice cloning.

Voice cloning, especially few-shot [280] and zero-shot TTS [78], [85], poses a significant threat to systems that rely on voice authentication, such as banking, customer service,

and other identity verification processes. With a convincing clone of someone's voice, attackers could potentially impersonate individuals to gain unauthorized access to sensitive information or accounts.

To address these concerns, it's essential to establish robust security protocols, consent-based regulations, and public awareness around voice cloning. Furthermore, advancements in detecting voice clones are equally important to help distinguish genuine voices from synthesized ones, protecting both individuals and organizations from potential misuse.

### B. Limitation of this Survey

Although we conduct a comprehensive survey of controllable TTS in this paper, there are some limitations we want to address in the future: 1) A unified benchmark method will be provided to evaluate controllable TTS methods. 2) Detailed analysis and control strategies of each specific speech attribute will be provided in an updated version of this paper. 3) The methodologies for feature disentanglement are crucial for controllable TTS but are not adequately discussed.

## VIII. CONCLUSION

In this survey paper, we first elaborate on the general pipeline for controllable TTS, followed by a glimpse of uncontrollable TTS methods. Then we comprehensively review existing controllable methods from the perspectives of model architectures and control strategies. Popular datasets and commonly used evaluation metrics for controllable TTS are also summarized in this paper. Besides, the current challenges are deeply analyzed and the promising future directions are also pointed out. To the best of our knowledge, this is the first comprehensive survey for controllable TTS.

Writing a comprehensive survey is a challenging task and an iterative process. Hence, we will regularly update this survey to offer researchers and practitioners a continuously evolving resource for understanding controllable speech synthesis.

## REFERENCES

[1] Wikipedia, "Speech synthesis." https://en.wikipedia.org/wiki/Speech_synthesis. Accessed: 2024-10-19.

[2] T. Dutoit, *An introduction to text-to-speech synthesis*, vol. 3. Springer Science & Business Media, 1997.

[3] S. Latif, J. Qadir, A. Qayyum, M. Usama, and S. Younis, "Speech technology for healthcare: Opportunities, challenges, and state of the art," *IEEE Reviews in Biomedical Engineering*, vol. 14, pp. 342–356, 2020.

[4] G. K. Anumanchipalli, J. Chartier, and E. F. Chang, "Speech synthesis from neural decoding of spoken sentences," *Nature*, vol. 568, no. 7753, pp. 493–498, 2019.

[5] G. López, L. Quesada, and L. A. Guerrero, "Alexa vs. siri vs. cortana vs. google assistant: a comparison of speech-based natural user interfaces," in *Advances in Human Factors and Systems Interaction: Proceedings of the AHFE 2017 International Conference on Human Factors and Systems Interaction*, pp. 241–250, 2018.

[6] Y. Li, F. Yu, Y.-Q. Xu, E. Chang, and H.-Y. Shum, "Speech-driven cartoon animation with emotions," in *Proceedings of the Ninth ACM International Conference on Multimedia*, pp. 365–371, 2001.

[7] Y. Wang, W. Wang, W. Liang, and L.-F. Yu, "Comic-guided speech synthesis," *ACM Transactions on Graphics*, vol. 38, no. 6, pp. 1–14, 2019.

[8] M. Marge, C. Espy-Wilson, N. G. Ward, A. Alwan, Y. Artzi, M. Bansal, G. Blankenship, J. Chai, H. Daumé III, D. Dey, *et al.*, "Spoken language interaction with robots: Recommendations for future research," *Computer Speech & Language*, vol. 71, p. 101255, 2022.

[9] S. Roehling, B. MacDonald, and C. Watson, "Towards expressive speech synthesis in English on a robotic platform," in *Proceedings of the Australasian International Conference on Speech Science and Technology*, pp. 130–135, 2006.

[10] OpenAI, "Introducing chatgpt." https://openai.com/index/chatgpt/. Accessed: 2024-10-22.

[11] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, *et al.*, "Llama: Open and efficient foundation language models," *arXiv preprint arXiv:2302.13971*, 2023.

[12] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.

[13] J. D. Owens, M. Houston, D. Luebke, S. Green, J. E. Stone, and J. C. Phillips, "Gpu computing," *Proceedings of the IEEE*, vol. 96, no. 5, pp. 879–899, 2008.

[14] X. Tan, J. Chen, H. Liu, J. Cong, C. Zhang, Y. Liu, X. Wang, Y. Leng, Y. Yi, L. He, *et al.*, "Naturalspeech: End-to-end text-to-speech synthesis with human-level quality," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.

[15] 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.

[16] S. Ö. Arık, 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 *International Conference on Machine Learning*, pp. 195–204, 2017.

[17] Z. Du, Q. Chen, S. Zhang, K. Hu, H. Lu, Y. Yang, H. Hu, S. Zheng, Y. Gu, Z. Ma, *et al.*, "Cosyvoice: A scalable multilingual zero-shot text-to-speech synthesizer based on supervised semantic tokens," *arXiv preprint arXiv:2407.05407*, 2024.

[18] Y.-J. Zhang, S. Pan, L. He, and Z.-H. Ling, "Learning latent representations for style control and transfer in end-to-end speech synthesis," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 6945–6949, 2019.

[19] 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*, pp. 5180–5189, 2018.

[20] S.-Y. Um, S. Oh, K. Byun, I. Jang, C. Ahn, and H.-G. Kang, "Emotional speech synthesis with rich and granularized control," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 7254–7258, 2020.

[21] T. Li, S. Yang, L. Xue, and L. Xie, "Controllable emotion transfer for end-to-end speech synthesis," in *12th International Symposium on Chinese Spoken Language Processing*, pp. 1–5, 2021.

[22] G. Zhang, Y. Qin, W. Zhang, J. Wu, M. Li, Y. Gai, F. Jiang, and T. Lee, "iemotts: Toward robust cross-speaker emotion transfer and control for speech synthesis based on disentanglement between prosody and timbre," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 31, pp. 1693–1705, 2023.

[23] J. Wang, Z. Wang, X. Hu, X. Li, Q. Fang, and L. Liu, "Residual-guided personalized speech synthesis based on face image," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 4743–4747, 2022.

[24] S. Goto, K. Onishi, Y. Saito, K. Tachibana, and K. Mori, "Face2speech: Towards multi-speaker text-to-speech synthesis using an embedding vector predicted from a face image.," in *INTERSPEECH*, pp. 1321–1325, 2020.

[25] J. Choi, J. Hong, and Y. M. Ro, "Diffv2s: Diffusion-based video-to-speech synthesis with vision-guided speaker embedding," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 7812–7821, 2023.

[26] A. Q. Jiang, A. Sablayrolles, A. Mensch, C. Bamford, D. S. Chaplot, D. d. l. Casas, F. Bressand, G. Lengyel, G. Lample, L. Saulnier, *et al.*, "Mistral 7b," *arXiv preprint arXiv:2310.06825*, 2023.

[27] J. Bai, S. Bai, Y. Chu, Z. Cui, K. Dang, X. Deng, Y. Fan, W. Ge, Y. Han, F. Huang, *et al.*, "Qwen technical report," *arXiv preprint arXiv:2309.16609*, 2023.

[28] X. Bi, D. Chen, G. Chen, S. Chen, D. Dai, C. Deng, H. Ding, K. Dong, Q. Du, Z. Fu, *et al.*, "Deepseek llm: Scaling open-source language models with longtermism," *arXiv preprint arXiv:2401.02954*, 2024.

[29] T. GLM, A. Zeng, B. Xu, B. Wang, C. Zhang, D. Yin, D. Rojas, G. Feng, H. Zhao, H. Lai, *et al.*, "Chatglm: A family of large language models from glm-130b to glm-4 all tools," *arXiv preprint arXiv:2406.12793*, 2024.

[30] H. Hao, L. Zhou, S. Liu, J. Li, S. Hu, R. Wang, and F. Wei, "Boosting large language model for speech synthesis: An empirical study," *arXiv preprint arXiv:2401.00246*, 2023.

[31] X. Gao, C. Zhang, Y. Chen, H. Zhang, and N. F. Chen, "Emodpo: Controllable emotional speech synthesis through direct preference optimization," *arXiv preprint arXiv:2409.10157*, 2024.

[32] P. Neekhara, S. Hussain, S. Ghosh, J. Li, R. Valle, R. Badlani, and B. Ginsburg, "Improving robustness of llm-based speech synthesis by learning monotonic alignment," *arXiv preprint arXiv:2406.17957*, 2024.

[33] Q. Fang, S. Guo, Y. Zhou, Z. Ma, S. Zhang, and Y. Feng, "Llama-omni: Seamless speech interaction with large language models," *arXiv preprint arXiv:2409.06666*, 2024.

[34] D. Zhang, S. Li, X. Zhang, J. Zhan, P. Wang, Y. Zhou, and X. Qiu, "Speechgpt: Empowering large language models with intrinsic cross-modal conversational abilities," *arXiv preprint arXiv:2305.11000*, 2023.

[35] X. Zhang, X. Lyu, Z. Du, Q. Chen, D. Zhang, H. Hu, C. Tan, T. Zhao, Y. Wang, B. Zhang, *et al.*, "Intrinsicvoice: Empowering llms with intrinsic real-time voice interaction abilities," *arXiv preprint arXiv:2410.08035*, 2024.

[36] D. H. Klatt, "Review of text-to-speech conversion for english," *The Journal of the Acoustical Society of America*, vol. 82, no. 3, pp. 737–793, 1987.

[37] T. Dutoit, "High-quality text-to-speech synthesis: An overview," *Journal Of Electrical And Electronics Engineering Australia*, vol. 17, no. 1, pp. 25–36, 1997.

[38] A. Breen, "Speech synthesis models: a review," *Electronics & Communication Engineering Journal*, vol. 4, no. 1, pp. 19–31, 1992.

[39] J. P. Olive and M. Y. Liberman, "Text to speech—an overview," *The Journal of the Acoustical Society of America*, vol. 78, no. S1, pp. S6–S6, 1985.

[40] S. King, "Measuring a decade of progress in text-to-speech," *Loquens*, vol. 1, no. 1, pp. e006–e006, 2014.

[41] H. Zen, K. Tokuda, and A. W. Black, "Statistical parametric speech synthesis," *Speech Communication*, vol. 51, no. 11, p. 1039 – 1064, 2009.

[42] X. Tan, T. Qin, F. Soong, and T.-Y. Liu, "A survey on neural speech synthesis," *arXiv preprint arXiv:2106.15561*, 2021.

[43] Y. Ning, S. He, Z. Wu, C. Xing, and L.-J. Zhang, "A review of deep learning based speech synthesis," *Applied Sciences*, vol. 9, no. 19, p. 4050, 2019.

[44] N. Kaur and P. Singh, "Conventional and contemporary approaches used in text to speech synthesis: A review," *Artificial Intelligence Review*, vol. 56, no. 7, pp. 5837–5880, 2023.

[45] W. Mattheyses and W. Verhelst, "Audiovisual speech synthesis: An overview of the state-of-the-art," *Speech Communication*, vol. 66, pp. 182–217, 2015.

[46] A. Triantafyllopoulos, B. W. Schuller, G. İymen, M. Sezgin, X. He, Z. Yang, P. Tzirakis, S. Liu, S. Mertes, E. André, *et al.*, "An overview of affective speech synthesis and conversion in the deep learning era," *Proceedings of the IEEE*, vol. 111, no. 10, pp. 1355–1381, 2023.

[47] Z. Mu, X. Yang, and Y. Dong, "Review of end-to-end speech synthesis technology based on deep learning," *arXiv preprint arXiv:2104.09995*, 2021.

[48] A. Mehrish, N. Majumder, R. Bharadwaj, R. Mihalcea, and S. Poria, "A review of deep learning techniques for speech processing," *Information Fusion*, vol. 99, p. 101869, 2023.

[49] Y. Tabet and M. Boughazi, "Speech synthesis techniques. a survey," in *International Workshop on Systems, Signal Processing and their Applications, WOSSPA*, pp. 67–70, 2011.

[50] C. Zhang, C. Zhang, S. Zheng, M. Zhang, M. Qamar, S.-H. Bae, and I. S. Kweon, "A survey on audio diffusion models: Text to speech synthesis and enhancement in generative ai," *arXiv preprint arXiv:2303.13336*, 2023.

[51] L. R. Rabiner, "Digital-formant synthesizer for speech-synthesis studies," *The Journal of the Acoustical Society of America*, vol. 43, no. 4, pp. 822–828, 1968.

[52] J. Allen, M. S. Hunnicutt, D. H. Klatt, R. C. Armstrong, and D. B. Pisoni, *From text to speech: The MITalk system*. Cambridge University Press, 1987.

[53] D. W. Purcell and K. G. Munhall, "Adaptive control of vowel formant frequency: Evidence from real-time formant manipulation," *The Journal of the Acoustical Society of America*, vol. 120, no. 2, pp. 966–977, 2006.

[54] D. H. Klatt, "Software for a cascade/parallel formant synthesizer," *Journal of the Acoustical Society of America*, vol. 67, no. 3, p. 971 – 995, 1980.

[55] J. Wouters and M. W. Macon, "Control of spectral dynamics in concatenative speech synthesis," *IEEE Transactions on Speech and Audio Processing*, vol. 9, no. 1, pp. 30–38, 2001.

[56] M. Bulut, S. S. Narayanan, and A. K. Syrdal, "Expressive speech synthesis using a concatenative synthesizer.," in *INTERSPEECH*, pp. 1265–1268, 2002.

[57] I. Bulyko and M. Ostendorf, "Joint prosody prediction and unit selection for concatenative speech synthesis," in *2001 IEEE International Conference on Acoustics, Speech, and Signal Processing. Proceedings (Cat. No. 01CH37221)*, vol. 2, pp. 781–784, 2001.

[58] Y. Stylianou, "Applying the harmonic plus noise model in concatenative speech synthesis," *IEEE Transactions on Speech and Audio Processing*, 2001.

[59] A. J. Hunt and A. W. Black, "Unit selection in a concatenative speech synthesis system using a large speech database," in *IEEE International Conference on Acoustics, Speech and Signal Processing - Proceedings*, 1996.

[60] T. Nose, J. Yamagishi, T. Masuko, and T. Kobayashi, "A style control technique for hmm-based expressive speech synthesis," *IEICE Transactions on Information and Systems*, vol. 90, no. 9, pp. 1406–1413, 2007.

[61] Z.-H. Ling, K. Richmond, J. Yamagishi, and R.-H. Wang, "Integrating articulatory features into hmm-based parametric speech synthesis," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 17, no. 6, pp. 1171–1185, 2009.

[62] T. Nose, M. Tachibana, and T. Kobayashi, "Hmm-based style control for expressive speech synthesis with arbitrary speaker's voice using model adaptation," *IEICE Transactions on Information and Systems*, vol. 92, no. 3, pp. 489–497, 2009.

[63] T. Yoshimura, K. Tokuda, T. Masuko, T. Kobayashi, and T. Kitamura, "Duration modeling for hmm-based speech synthesis.," in *ICSLP*, vol. 98, pp. 29–32, 1998.

[64] T. Yoshimura, K. Tokuda, T. Masuko, T. Kobayashi, and T. Kitamura, "Simultaneous modeling of spectrum, pitch and duration in hmm-based speech synthesis," in *Sixth European conference on speech communication and technology*, 1999.

[65] K. Tokuda, T. Yoshimura, T. Masuko, T. Kobayashi, and T. Kitamura, "Speech parameter generation algorithms for hmm-based speech synthesis," in *2000 IEEE international conference on acoustics, speech, and signal processing. Proceedings (Cat. No. 00CH37100)*, vol. 3, pp. 1315–1318, IEEE, 2000.

[66] J. Yamagishi, T. Nose, H. Zen, Z.-H. Ling, T. Toda, K. Tokuda, S. King, and S. Renals, "Robust speaker-adaptive hmm-based text-to-speech synthesis," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 17, no. 6, pp. 1208–1230, 2009.

[67] J. Yamagishi, T. Kobayashi, Y. Nakano, K. Ogata, and J. Isogai, "Analysis of speaker adaptation algorithms for hmm-based speech synthesis and a constrained smaplr adaptation algorithm," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 17, no. 1, pp. 66–83, 2009.

[68] T. Masuko, K. Tokuda, T. Kobayashi, and S. Imai, "Voice characteristics conversion for hmm-based speech synthesis system," in *1997 IEEE international conference on acoustics, speech, and signal processing*, vol. 3, pp. 1611–1614, 1997.

[69] C.-H. Wu, C.-C. Hsia, T.-H. Liu, and J.-F. Wang, "Voice conversion using duration-embedded bi-hmms for expressive speech synthesis," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 14, no. 4, pp. 1109–1116, 2006.

[70] J. Yamagishi, K. Onishi, T. Masuko, and T. Kobayashi, "Modeling of various speaking styles and emotions for hmm-based speech synthesis.," in *interspeech*, pp. 2461–2464, 2003.

[71] J. Yamagishi, K. Onishi, T. Masuko, and T. Kobayashi, "Acoustic modeling of speaking styles and emotional expressions in hmm-based speech synthesis," *IEICE Transactions on Information and Systems*, vol. 88, no. 3, pp. 502–509, 2005.

[72] J. Lorenzo-Trueba, R. Barra-Chicote, R. San-Segundo, J. Ferreiros, J. Yamagishi, and J. M. Montero, "Emotion transplantation through adaptation in hmm-based speech synthesis," *Computer Speech & Language*, vol. 34, no. 1, pp. 292–307, 2015.

[73] A. Van Den Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, K. Kavukcuoglu, *et al.*, "Wavenet: A generative model for raw audio," *arXiv preprint arXiv:1609.03499*, vol. 12, 2016.

[74] 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.

[75] J. Shen, R. Pang, R. J. Weiss, M. Schuster, N. Jaitly, Z. Yang, Z. Chen, Y. Zhang, Y. Wang, R. Skerrv-Ryan, *et al.*, "Natural tts synthesis by conditioning wavenet on mel spectrogram predictions," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 4779–4783, 2018.

[76] 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.

[77] A. Łańcucki, "Fastpitch: Parallel text-to-speech with pitch prediction," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 6588–6592, 2021.

[78] Y. Wang, H. Zhan, L. Liu, R. Zeng, H. Guo, J. Zheng, Q. Zhang, X. Zhang, S. Zhang, and Z. Wu, "Maskgct: Zero-shot text-to-speech with masked generative codec transformer," *arXiv preprint arXiv:2409.00750*, 2024.

[79] Y. Fan, Y. Qian, F. K. Soong, and L. He, "Multi-speaker modeling and speaker adaptation for dnn-based tts synthesis," in *IEEE international conference on acoustics, speech and signal processing*, pp. 4475–4479, 2015.

[80] S.-F. Huang, C.-J. Lin, D.-R. Liu, Y.-C. Chen, and H.-y. Lee, "Meta-tts: Meta-learning for few-shot speaker adaptive text-to-speech," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 30, pp. 1558–1571, 2022.

[81] M. Chen, X. Tan, Y. Ren, J. Xu, H. Sun, S. Zhao, T. Qin, and T.-Y. Liu, "Multispeech: Multi-speaker text to speech with transformer," *arXiv preprint arXiv:2006.04664*, 2020.

[82] E. Casanova, J. Weber, C. D. Shulby, A. C. Junior, E. Gölge, and M. A. Ponti, "Yourtts: Towards zero-shot multi-speaker tts and zero-shot voice conversion for everyone," in *International Conference on Machine Learning*, pp. 2709–2720, 2022.

[83] Y. Lei, S. Yang, X. Wang, and L. Xie, "Msemotts: Multi-scale emotion transfer, prediction, and control for emotional speech synthesis," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 30, pp. 853–864, 2022.

[84] I. Elias, H. Zen, J. Shen, Y. Zhang, Y. Jia, R. J. Weiss, and Y. Wu, "Parallel tacotron: Non-autoregressive and controllable tts," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 5709–5713, 2021.

[85] C. Wang, S. Chen, Y. Wu, Z. Zhang, L. Zhou, S. Liu, Z. Chen, Y. Liu, H. Wang, J. Li, *et al.*, "Neural codec language models are zero-shot text to speech synthesizers," *arXiv preprint arXiv:2301.02111*, 2023.

[86] K. Shen, Z. Ju, X. Tan, Y. Liu, Y. Leng, L. He, T. Qin, S. Zhao, and J. Bian, "Naturalspeech 2: Latent diffusion models are natural and zero-shot speech and singing synthesizers," *arXiv preprint arXiv:2304.09116*, 2023.

[87] Z. Ju, Y. Wang, K. Shen, X. Tan, D. Xin, D. Yang, Y. Liu, Y. Leng, K. Song, S. Tang, *et al.*, "Naturalspeech 3: Zero-shot speech synthesis with factorized codec and diffusion models," *arXiv preprint arXiv:2403.03100*, 2024.

[88] Y. A. Li, C. Han, and N. Mesgarani, "Styletts: A style-based generative model for natural and diverse text-to-speech synthesis," *arXiv preprint arXiv:2205.15439*, 2022.

[89] Y. A. Li, C. Han, V. Raghavan, G. Mischler, and N. Mesgarani, "Styletts 2: Towards human-level text-to-speech through style diffusion and adversarial training with large speech language models," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[90] R. Huang, Y. Ren, J. Liu, C. Cui, and Z. Zhao, "Generspeech: Towards style transfer for generalizable out-of-domain text-to-speech," *Advances in Neural Information Processing Systems*, vol. 35, pp. 10970–10983, 2022.

[91] Z. Jiang, Y. Ren, Z. Ye, J. Liu, C. Zhang, Q. Yang, S. Ji, R. Huang, C. Wang, X. Yin, *et al.*, "Mega-tts: Zero-shot text-to-speech at scale with intrinsic inductive bias," *arXiv preprint arXiv:2306.03509*, 2023.

[92] E. Cooper, C.-I. Lai, Y. Yasuda, F. Fang, X. Wang, N. Chen, and J. Yamagishi, "Zero-shot multi-speaker text-to-speech with state-of-the-art neural speaker embeddings," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 6184–6188, 2020.

[93] P. Peng, P.-Y. Huang, S.-W. Li, A. Mohamed, and D. Harwath, "Voicecraft: Zero-shot speech editing and text-to-speech in the wild," *arXiv preprint arXiv:2403.16973*, 2024.

[94] D. Tan, L. Deng, Y. T. Yeung, X. Jiang, X. Chen, and T. Lee, "Editspeech: A text based speech editing system using partial inference and bidirectional fusion," in *IEEE Automatic Speech Recognition and Understanding Workshop*, pp. 626–633, 2021.

[95] J. Tae, H. Kim, and T. Kim, "Editts: Score-based editing for controllable text-to-speech," *arXiv preprint arXiv:2110.02584*, 2021.

[96] S. Seshadri, T. Raitio, D. Castellani, and J. Li, "Emphasis control for parallel neural tts," *arXiv preprint arXiv:2110.03012*, 2021.

[97] T. B. Brown, "Language models are few-shot learners," *arXiv preprint arXiv:2005.14165*, 2020.

[98] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat, *et al.*, "Gpt-4 technical report," *arXiv preprint arXiv:2303.08774*, 2023.

[99] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, "Exploring the limits of transfer learning with a unified text-to-text transformer," *Journal of Machine Learning Research*, vol. 21, no. 140, pp. 1–67, 2020.

[100] A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, P. Barham, H. W. Chung, C. Sutton, S. Gehrmann, *et al.*, "Palm: Scaling language modeling with pathways," *Journal of Machine Learning Research*, vol. 24, no. 240, pp. 1–113, 2023.

[101] Z. Guo, Y. Leng, Y. Wu, S. Zhao, and X. Tan, "Promptts: Controllable text-to-speech with text descriptions," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 1–5, 2023.

[102] Y. Leng, Z. Guo, K. Shen, X. Tan, Z. Ju, Y. Liu, Y. Liu, D. Yang, L. Zhang, K. Song, *et al.*, "Promptts 2: Describing and generating voices with text prompt," *arXiv preprint arXiv:2309.02285*, 2023.

[103] Y. Zhou, X. Qin, Z. Jin, S. Zhou, S. Lei, S. Zhou, Z. Wu, and J. Jia, "Voxinstruct: Expressive human instruction-to-speech generation with unified multilingual codec language modelling," in *Proceedings of the 32nd ACM International Conference on Multimedia*, pp. 554–563, 2024.

[104] R. Shimizu, R. Yamamoto, M. Kawamura, Y. Shirahata, H. Doi, T. Komatsu, and K. Tachibana, "Promptts++: Controlling speaker identity in prompt-based text-to-speech using natural language descriptions," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 12672–12676, 2024.

[105] D. Yang, S. Liu, R. Huang, C. Weng, and H. Meng, "Instructtts: Modelling expressive tts in discrete latent space with natural language style prompt," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2024.

[106] S. Ji, J. Zuo, W. Wang, M. Fang, S. Zheng, Q. Chen, Z. Jiang, H. Huang, Z. Wang, X. Cheng, *et al.*, "Controlspeech: Towards simultaneous zero-shot speaker cloning and zero-shot language style control with decoupled codec," *arXiv preprint arXiv:2406.01205*, 2024.

[107] T. Fukada, K. Tokuda, T. Kobayashi, and S. Imai, "An adaptive algorithm for mel-cepstral analysis of speech.," in *icassp*, vol. 92, pp. 137–140, 1992.

[108] F. Itakura, "Line spectrum representation of linear predictor coefficients of speech signals," *The Journal of the Acoustical Society of America*, vol. 57, no. S1, pp. S35–S35, 1975.

[109] H. Kawahara, I. Masuda-Katsuse, and A. De Cheveigne, "Restructuring speech representations using a pitch-adaptive time–frequency smoothing and an instantaneous-frequency-based f0 extraction: Possible role of a repetitive structure in sounds," *Speech Communication*, vol. 27, no. 3-4, pp. 187–207, 1999.

[110] Wikipedia, "Spectrogram." https://en.wikipedia.org/wiki/Spectrogram. Accessed: 2024-10-24.

[111] H. Zen, A. Senior, and M. Schuster, "Statistical parametric speech synthesis using deep neural networks," in *2013 ieee international conference on acoustics, speech and signal processing*, pp. 7962–7966, IEEE, 2013.

[112] Y. Fan, Y. Qian, F.-L. Xie, and F. K. Soong, "Tts synthesis with bidirectional lstm based recurrent neural networks," in *Fifteenth annual conference of the international speech communication association*, 2014.

[113] M. Jeong, H. Kim, S. J. Cheon, B. J. Choi, and N. S. Kim, "Difftts: A denoising diffusion model for text-to-speech," *arXiv preprint arXiv:2104.01409*, 2021.

[114] X. Chen, L. Xu, Z. Liu, M. Sun, and H. Luan, "Joint learning of character and word embeddings," in *Twenty-fourth International Joint Conference on Artificial Intelligence*, 2015.

[115] F. Almeida and G. Xexéo, "Word embeddings: A survey," *arXiv preprint arXiv:1901.09069*, 2019.

[116] 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. 17022–17033, 2020.

[117] K. Tokuda, T. Yoshimura, T. Masuko, T. Kobayashi, and T. Kitamura, "Speech parameter generation algorithms for hmm-based speech synthesis," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, vol. 3, p. 1315 – 1318, 2000.

[118] H. Zen, T. Nose, J. Yamagishi, S. Sako, T. Masuko, A. W. Black, and K. Tokuda, "The hmm-based speech synthesis system (hts) version 2.0," in *6th ISCA Workshop on Speech Synthesis*, p. 294 – 299, 2007.

[119] T. Nose and T. Kobayashi, "An intuitive style control technique in hmm-based expressive speech synthesis using subjective style intensity and multiple-regression global variance model," *Speech Communication*, vol. 55, no. 2, pp. 347–357, 2013.

[120] Y. Nishigaki, S. Takamichi, T. Toda, G. Neubig, S. Sakti, and S. Nakamura, "Prosody-controllable hmm-based speech synthesis using speech input," *Proc. MLSLP*, 2015.

[121] S. Vasquez and M. Lewis, "Melnet: A generative model for audio in the frequency domain," *arXiv preprint arXiv:1906.01083*, 2019.

[122] S. Ö. Arık, 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 *International conference on machine learning*, pp. 195–204, 2017.

[123] K. Peng, W. Ping, Z. Song, and K. Zhao, "Non-autoregressive neural text-to-speech," in *International conference on machine learning*, pp. 7586–7598, 2020.

[124] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. u. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems*, vol. 30, pp. 1–11, 2017.

[125] 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*, pp. 6706–6713, 2019.

[126] J. Devlin, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.

[127] S.-H. Lee, H.-W. Yoon, H.-R. Noh, J.-H. Kim, and S.-W. Lee, "Multi-spectrogan: High-diversity and high-fidelity spectrogram generation with adversarial style combination for speech synthesis," in *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 13198–13206, 2021.

[128] S. Ma, D. Mcduff, and Y. Song, "Neural tts stylization with adversarial and collaborative games," in *International conference on learning representations*, 2018.

[129] H. Guo, F. K. Soong, L. He, and L. Xie, "A new gan-based end-to-end tts training algorithm," *arXiv preprint arXiv:1904.04775*, 2019.

[130] 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.

[131] V. Popov, I. Vovk, V. Gogoryan, T. Sadekova, and M. Kudinov, "Grad-tts: A diffusion probabilistic model for text-to-speech," in *International Conference on Machine Learning*, pp. 8599–8608, 2021.

[132] C. Miao, S. Liang, M. Chen, J. Ma, S. Wang, and J. Xiao, "Flow-tts: A non-autoregressive network for text to speech based on flow," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 7209–7213, IEEE, 2020.

[133] J. Kim, S. Kim, J. Kong, and S. Yoon, "Glow-tts: A generative flow for text-to-speech via monotonic alignment search," *Advances in Neural Information Processing Systems*, vol. 33, pp. 8067–8077, 2020.

[134] H. Kawahara, "Straight, exploitation of the other aspect of vocoder: Perceptually isomorphic decomposition of speech sounds," *Acoustical Science and Technology*, vol. 27, no. 6, pp. 349–353, 2006.

[135] 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.

[136] S. Mehri, K. Kumar, I. Gulrajani, R. Kumar, S. Jain, J. Sotelo, A. Courville, and Y. Bengio, "Samplernn: An unconditional end-to-end neural audio generation model," *arXiv preprint arXiv:1612.07837*, 2016.

[137] N. Kalchbrenner, E. Elsen, K. Simonyan, S. Noury, N. Casagrande, E. Lockhart, F. Stimberg, A. Oord, S. Dieleman, and K. Kavukcuoglu, "Efficient neural audio synthesis," in *International Conference on Machine Learning*, pp. 2410–2419, 2018.

[138] J.-M. Valin and J. Skoglund, "Lpcnet: Improving neural speech synthesis through linear prediction," in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 5891–5895, 2019.

[139] C. Yu, H. Lu, N. Hu, M. Yu, C. Weng, K. Xu, P. Liu, D. Tuo, S. Kang, G. Lei, *et al.*, "Durian: Duration informed attention network for speech synthesis.," in *Interspeech*, pp. 2027–2031, 2020.

[140] A. Oord, Y. Li, I. Babuschkin, K. Simonyan, O. Vinyals, K. Kavukcuoglu, G. Driessche, E. Lockhart, L. Cobo, F. Stimberg, *et al.*, "Parallel wavenet: Fast high-fidelity speech synthesis," in *International Conference on Machine Learning*, pp. 3918–3926, 2018.

[141] Z. Jin, A. Finkelstein, G. J. Mysore, and J. Lu, "Fftnet: A real-time speaker-dependent neural vocoder," in *2018 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pp. 2251–2255, 2018.

[142] C. Donahue, J. McAuley, and M. Puckette, "Adversarial audio synthesis," *arXiv preprint arXiv:1802.04208*, 2018.

[143] M. Bińkowski, J. Donahue, S. Dieleman, A. Clark, E. Elsen, N. Casagrande, L. C. Cobo, and K. Simonyan, "High fidelity speech synthesis with adversarial networks," *arXiv preprint arXiv:1909.11646*, 2019.

[144] R. Yamamoto, E. Song, and J.-M. Kim, "Parallel wavegan: A fast waveform generation model based on generative adversarial networks with multi-resolution spectrogram," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 6199–6203, 2020.

[145] 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.

[146] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," *Advances in Neural Information Processing Systems*, vol. 33, pp. 6840–6851, 2020.

[147] R. Huang, M. W. Lam, J. Wang, D. Su, D. Yu, Y. Ren, and Z. Zhao, "Fastdiff: A fast conditional diffusion model for high-quality speech synthesis," *arXiv preprint arXiv:2204.09934*, 2022.

[148] Z. Kong, W. Ping, J. Huang, K. Zhao, and B. Catanzaro, "Diffwave: A versatile diffusion model for audio synthesis," *arXiv preprint arXiv:2009.09761*, 2020.

[149] N. Chen, Y. Zhang, H. Zen, R. J. Weiss, M. Norouzi, and W. Chan, "Wavegrad: Estimating gradients for waveform generation," *arXiv preprint arXiv:2009.00713*, 2020.

[150] S.-g. Lee, H. Kim, C. Shin, X. Tan, C. Liu, Q. Meng, T. Qin, W. Chen, S. Yoon, and T.-Y. Liu, "Priorgrad: Improving conditional denoising diffusion models with data-dependent adaptive prior," *arXiv preprint arXiv:2106.06406*, 2021.

[151] Y. Ren, X. Tan, T. Qin, Z. Zhao, and T.-Y. Liu, "Revisiting over-smoothness in text to speech," *arXiv preprint arXiv:2202.13066*, 2022.

[152] S. Kim, K. Shih, J. F. Santos, E. Bakhturina, M. Desta, R. Valle, S. Yoon, B. Catanzaro, *et al.*, "P-flow: a fast and data-efficient zero-shot tts through speech prompting," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[153] Y. Guo, C. Du, Z. Ma, X. Chen, and K. Yu, "Voiceflow: Efficient text-to-speech with rectified flow matching," in *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 11121–11125, 2024.

[154] S.-H. Lee, H.-Y. Choi, and S.-W. Lee, "Periodwave: Multi-period flow matching for high-fidelity waveform generation," *arXiv preprint arXiv:2408.07547*, 2024.

[155] W. Ping, K. Peng, K. Zhao, and Z. Song, "Waveflow: A compact flow-based model for raw audio," in *International Conference on Machine Learning*, pp. 7706–7716, 2020.

[156] S. Kim, S.-g. Lee, J. Song, J. Kim, and S. Yoon, "Flowavenet: A generative flow for raw audio," *arXiv preprint arXiv:1811.02155*, 2018.

[157] K. Peng, W. Ping, Z. Song, and K. Zhao, "Non-autoregressive neural text-to-speech," in *International conference on machine learning*, pp. 7586–7598, 2020.

[158] H. Guo, F. Xie, X. Wu, F. K. Soong, and H. Meng, "Msmc-tts: Multi-stage multi-codebook vq-vae based neural tts," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 31, pp. 1811–1824, 2023.

[159] J. Kim, J. Kong, and J. Son, "Conditional variational autoencoder with adversarial learning for end-to-end text-to-speech," in *International Conference on Machine Learning*, pp. 5530–5540, 2021.

[160] J. Sotelo, S. Mehri, K. Kumar, J. F. Santos, K. Kastner, A. C. Courville, and Y. Bengio, "Char2wav: End-to-end speech synthesis," in *5th International Conference on Learning Representations, Workshop Track Proceedings*, 2017.

[161] W. Ping, K. Peng, and J. Chen, "Clarinet: Parallel wave generation in end-to-end text-to-speech," *arXiv preprint arXiv:1807.07281*, 2018.

[162] J. Donahue, S. Dieleman, M. Bińkowski, E. Elsen, and K. Simonyan, "End-to-end adversarial text-to-speech," *arXiv preprint arXiv:2006.03575*, 2020.

[163] D. Rezende and S. Mohamed, "Variational inference with normalizing flows," in *International conference on machine learning*, pp. 1530–1538, 2015.

[164] A. Baevski, S. Schneider, and M. Auli, "vq-wav2vec: Self-supervised learning of discrete speech representations," *arXiv preprint arXiv:1910.05453*, 2019.

[165] A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," *Advances in neural information processing systems*, vol. 33, pp. 12449–12460, 2020.

[166] W.-N. Hsu, B. Bolte, Y.-H. H. Tsai, K. Lakhotia, 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.

[167] L. Barrault, Y.-A. Chung, M. C. Meglioli, D. Dale, N. Dong, M. Duppenthaler, P.-A. Duquenne, B. Ellis, H. Elsahar, H. Haaheim, *et al.*, "Seamless: Multilingual expressive and streaming speech translation," *arXiv preprint arXiv:2312.05187*, 2023.

[168] N. Zeghidour, A. Luebs, A. Omran, J. Skoglund, and M. Tagliasacchi, "Soundstream: An end-to-end neural audio codec," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 30, pp. 495–507, 2021.

[169] A. Défossez, J. Copet, G. Synnaeve, and Y. Adi, "High fidelity neural audio compression," *arXiv preprint arXiv:2210.13438*, 2022.

[170] D. Yang, S. Liu, R. Huang, J. Tian, C. Weng, and Y. Zou, "Hifi-codec: Group-residual vector quantization for high fidelity audio codec," *arXiv preprint arXiv:2305.02765*, 2023.

[171] X. Zhang, D. Zhang, S. Li, Y. Zhou, and X. Qiu, "Speechtokenizer: Unified speech tokenizer for speech large language models," *arXiv preprint arXiv:2308.16692*, 2023.

[172] R. Kumar, P. Seetharaman, A. Luebs, I. Kumar, and K. Kumar, "High-fidelity audio compression with improved rvqgan," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[173] K. Tokuda, Y. Nankaku, T. Toda, H. Zen, J. Yamagishi, and K. Oura, "Speech synthesis based on hidden markov models," *Proceedings of the IEEE*, vol. 101, no. 5, pp. 1234–1252, 2013.

[174] H. Zen, K. Tokuda, and A. W. Black, "Statistical parametric speech synthesis," *speech communication*, vol. 51, no. 11, pp. 1039–1064, 2009.

[175] J. Shen, R. Pang, R. J. Weiss, M. Schuster, N. Jaitly, Z. Yang, Z. Chen, Y. Zhang, Y. Wang, R. Skerrv-Ryan, *et al.*, "Natural tts synthesis by conditioning wavenet on mel spectrogram predictions," in *2018 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pp. 4779–4783, IEEE, 2018.

[176] 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," *arXiv preprint arXiv:1710.07654*, 2017.

[177] 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*, pp. 6706–6713, 2019.

[178] N. Chen, Y. Zhang, H. Zen, R. J. Weiss, M. Norouzi, N. Dehak, and W. Chan, "Wavegrad 2: Iterative refinement for text-to-speech synthesis," *arXiv preprint arXiv:2106.09660*, 2021.

[179] L. Qin, Z.-H. Ling, Y.-J. Wu, B.-F. Zhang, and R.-H. Wang, "Hmm-based emotional speech synthesis using average emotion model," in *Chinese Spoken Language Processing: 5th International Symposium, ISCSLP 2006, Singapore, December 13-16, 2006. Proceedings*, pp. 233–240, Springer, 2006.

[180] 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.

[181] Z. Liu, Q. Tian, C. Hu, X. Liu, M. Wu, Y. Wang, H. Zhao, and Y. Wang, "Controllable and lossless non-autoregressive end-to-end text-to-speech," *arXiv preprint arXiv:2207.06088*, 2022.

[182] M. Kang, D. Min, and S. J. Hwang, "Grad-stylespeech: Any-speaker adaptive text-to-speech synthesis with diffusion models," in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1–5, IEEE, 2023.

[183] H. J. Park, J. S. Kim, W. Shin, and S. W. Han, "Dex-tts: Diffusion-based expressive text-to-speech with style modeling on time variability," *arXiv preprint arXiv:2406.19135*, 2024.

[184] Y. Gao, N. Morioka, Y. Zhang, and N. Chen, "E3 tts: Easy end-to-end diffusion-based text to speech," in *2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pp. 1–8, IEEE, 2023.

[185] K. Lee, D. W. Kim, J. Kim, and J. Cho, "Ditto-tts: Efficient and scalable zero-shot text-to-speech with diffusion transformer," *arXiv preprint arXiv:2406.11427*, 2024.

[186] D. Yang, D. Wang, H. Guo, X. Chen, X. Wu, and H. Meng, "Simplespeech: Towards simple and efficient text-to-speech with scalar latent transformer diffusion models," *arXiv preprint arXiv:2406.02328*, 2024.

[187] A. Vyas, B. Shi, M. Le, A. Tjandra, Y.-C. Wu, B. Guo, J. Zhang, X. Zhang, R. Adkins, W. Ngan, *et al.*, "Audiobox: Unified audio generation with natural language prompts," *arXiv preprint arXiv:2312.15821*, 2023.

[188] S. Kim, K. Shih, J. F. Santos, E. Bakhturina, M. Desta, R. Valle, S. Yoon, B. Catanzaro, *et al.*, "P-flow: a fast and data-efficient zero-shot tts through speech prompting," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[189] M. Le, A. Vyas, B. Shi, B. Karrer, L. Sari, R. Moritz, M. Williamson, V. Manohar, Y. Adi, J. Mahadeokar, *et al.*, "Voicebox: Text-guided multilingual universal speech generation at scale," *Advances in neural information processing systems*, vol. 36, 2024.

[190] A. H. Liu, M. Le, A. Vyas, B. Shi, A. Tjandra, and W.-N. Hsu, "Generative pre-training for speech with flow matching," *arXiv preprint arXiv:2310.16338*, 2023.

[191] S.-H. Lee, H.-Y. Choi, S.-B. Kim, and S.-W. Lee, "Hierspeech++: Bridging the gap between semantic and acoustic representation of speech by hierarchical variational inference for zero-shot speech synthesis," *arXiv preprint arXiv:2311.12454*, 2023.

[192] Z. Ye, Z. Ju, H. Liu, X. Tan, J. Chen, Y. Lu, P. Sun, J. Pan, W. Bian, S. He, *et al.*, "Flashspeech: Efficient zero-shot speech synthesis," in *Proceedings of the 32nd ACM International Conference on Multimedia*, pp. 6998–7007, 2024.

[193] S. E. Eskimez, X. Wang, M. Thakker, C. Li, C.-H. Tsai, Z. Xiao, H. Yang, Z. Zhu, M. Tang, X. Tan, *et al.*, "E2 tts: Embarrassingly easy fully non-autoregressive zero-shot tts," *arXiv preprint arXiv:2406.18009*, 2024.

[194] Y. Chen, Z. Niu, Z. Ma, K. Deng, C. Wang, J. Zhao, K. Yu, and X. Chen, "F5-tts: A fairytaler that fakes fluent and faithful speech with flow matching," *arXiv preprint arXiv:2410.06885*, 2024.

[195] S. Woo, S. Debnath, R. Hu, X. Chen, Z. Liu, I. S. Kweon, and S. Xie, "Convnext v2: Co-designing and scaling convnets with masked autoencoders," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16133–16142, 2023.

[196] Z. Liu, S. Wang, P. Zhu, M. Bi, and H. Li, "E1 tts: Simple and fast non-autoregressive tts," *arXiv preprint arXiv:2409.09351*, 2024.

[197] W. Luo, T. Hu, S. Zhang, J. Sun, Z. Li, and Z. Zhang, "Diff-instruct: A universal approach for transferring knowledge from pre-trained diffusion models," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[198] T. Yin, M. Gharbi, T. Park, R. Zhang, E. Shechtman, F. Durand, and W. T. Freeman, "Improved distribution matching distillation for fast image synthesis," *arXiv preprint arXiv:2405.14867*, 2024.

[199] D. Yang, R. Huang, Y. Wang, H. Guo, D. Chong, S. Liu, X. Wu, and H. Meng, "Simplespeech 2: Towards simple and efficient text-to-speech with flow-based scalar latent transformer diffusion models," *arXiv preprint arXiv:2408.13893*, 2024.

[200] S. Chen, C. Wang, Z. Chen, Y. Wu, S. Liu, Z. Chen, J. Li, N. Kanda, T. Yoshioka, X. Xiao, *et al.*, "Wavlm: Large-scale self-supervised pre-training for full stack speech processing," *IEEE Journal of Selected Topics in Signal Processing*, vol. 16, no. 6, pp. 1505–1518, 2022.

[201] 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*, pp. 4693–4702, PMLR, 2018.

[202] A. Défossez, J. Copet, G. Synnaeve, and Y. Adi, "High fidelity neural audio compression," *arXiv preprint arXiv:2210.13438*, 2022.

[203] E. Kharitonov, D. Vincent, Z. Borsos, R. Marinier, S. Girgin, O. Pietquin, M. Sharifi, M. Tagliasacchi, and N. Zeghidour, "Speak, read and prompt: High-fidelity text-to-speech with minimal supervision," *Transactions of the Association for Computational Linguistics*, vol. 11, pp. 1703–1718, 2023.

[204] R. Huang, C. Zhang, Y. Wang, D. Yang, L. Liu, Z. Ye, Z. Jiang, C. Weng, Z. Zhao, and D. Yu, "Make-a-voice: Unified voice synthesis with discrete representation," *arXiv preprint arXiv:2305.19269*, 2023.

[205] D. Kim, S. Hong, and Y.-H. Choi, "Sc vall-e: Style-controllable zero-shot text to speech synthesizer," *arXiv preprint arXiv:2307.10550*, 2023.

[206] S. Chen, Y. Feng, L. He, T. He, W. He, Y. Hu, B. Lin, Y. Lin, Y. Pan, P. Tan, *et al.*, "Takin: A cohort of superior quality zero-shot speech generation models," *arXiv preprint arXiv:2409.12139*, 2024.

[207] Z. Zhang, L. Zhou, C. Wang, S. Chen, Y. Wu, S. Liu, Z. Chen, Y. Liu, H. Wang, J. Li, *et al.*, "Speak foreign languages with your own voice: Cross-lingual neural codec language modeling," *arXiv preprint arXiv:2303.03926*, 2023.

[208] Y. Song, Z. Chen, X. Wang, Z. Ma, and X. Chen, "Ella-v: Stable neural codec language modeling with alignment-guided sequence reordering," *arXiv preprint arXiv:2401.07333*, 2024.

[209] D. Xin, X. Tan, K. Shen, Z. Ju, D. Yang, Y. Wang, S. Takamichi, H. Saruwatari, S. Liu, J. Li, *et al.*, "Rall-e: Robust codec language modeling with chain-of-thought prompting for text-to-speech synthesis," *arXiv preprint arXiv:2404.03204*, 2024.

[210] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou, *et al.*, "Chain-of-thought prompting elicits reasoning in large language models," *Advances in neural information processing systems*, vol. 35, pp. 24824–24837, 2022.

[211] B. Han, L. Zhou, S. Liu, S. Chen, L. Meng, Y. Qian, Y. Liu, S. Zhao, J. Li, and F. Wei, "Vall-e r: Robust and efficient zero-shot text-to-speech synthesis via monotonic alignment," *arXiv preprint arXiv:2406.07855*, 2024.

[212] S. Chen, S. Liu, L. Zhou, Y. Liu, X. Tan, J. Li, S. Zhao, Y. Qian, and F. Wei, "Vall-e 2: Neural codec language models are human parity zero-shot text to speech synthesizers," *arXiv preprint arXiv:2406.05370*, 2024.

[213] Y. Nishimura, T. Hirose, M. Ohi, H. Nakayama, and N. Inoue, "Hall-e: Hierarchical neural codec language model for minute-long zero-shot text-to-speech synthesis," *arXiv preprint arXiv:2410.04380*, 2024.

[214] H.-H. Guo, K. Liu, F.-Y. Shen, Y.-C. Wu, F.-L. Xie, K. Xie, and K.-T. Xu, "Fireredtts: A foundation text-to-speech framework for industry-level generative speech applications," *arXiv preprint arXiv:2409.03283*, 2024.

[215] H. Guo, F. Xie, D. Yang, X. Wu, and H. Meng, "Speaking from coarse to fine: Improving neural codec language model via multi-scale speech coding and generation," *arXiv preprint arXiv:2409.11630*, 2024.

[216] M. Łajszczak, G. Cámbara, Y. Li, F. Beyhan, A. van Korlaar, F. Yang, A. Joly, Á. Martín-Cortinas, A. Abbas, A. Michalski, *et al.*, "Base tts: Lessons from building a billion-parameter text-to-speech model on 100k hours of data," *arXiv preprint arXiv:2402.08093*, 2024.

[217] P. Anastassiou, J. Chen, J. Chen, Y. Chen, Z. Chen, Z. Chen, J. Cong, L. Deng, C. Ding, L. Gao, *et al.*, "Seed-tts: A family of high-quality versatile speech generation models," *arXiv preprint arXiv:2406.02430*, 2024.

[218] X. Di, Z. Chen, Y. Liang, J. Zheng, Y. Wang, and C. Ding, "Bailing-tts: Chinese dialectal speech synthesis towards human-like spontaneous representation," *arXiv preprint arXiv:2408.00284*, 2024.

[219] L. Meng, L. Zhou, S. Liu, S. Chen, B. Han, S. Hu, Y. Liu, J. Li, S. Zhao, X. Wu, *et al.*, "Autoregressive speech synthesis without vector quantization," *arXiv preprint arXiv:2407.08551*, 2024.

[220] Z. Liu, S. Wang, S. Inoue, Q. Bai, and H. Li, "Autoregressive diffusion transformer for text-to-speech synthesis," *arXiv preprint arXiv:2406.05551*, 2024.

[221] S. Vekkot, D. Gupta, M. Zakariah, and Y. A. Alotaibi, "Emotional voice conversion using a hybrid framework with speaker-adaptive dnn and particle-swarm-optimized neural network," *IEEE Access*, vol. 8, pp. 74627–74647, 2020.

[222] I. Elias, H. Zen, J. Shen, Y. Zhang, Y. Jia, R. J. Weiss, and Y. Wu, "Parallel tacotron: Non-autoregressive and controllable tts," in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 5709–5713, IEEE, 2021.

[223] M. Kim, S. J. Cheon, B. J. Choi, J. J. Kim, and N. S. Kim, "Expressive text-to-speech using style tag," *INTERSPEECH 2021*, pp. 4663–4667, 2021.

[224] E. Casanova, C. Shulby, E. Gölge, N. M. Müller, F. S. De Oliveira, A. C. Junior, A. d. S. Soares, S. M. Aluisio, and M. A. Ponti, "Sc-glowtts: An efficient zero-shot multi-speaker text-to-speech model," *arXiv preprint arXiv:2104.05557*, 2021.

[225] D. Min, D. B. Lee, E. Yang, and S. J. Hwang, "Meta-stylespeech: Multi-speaker adaptive text-to-speech generation," in *International Conference on Machine Learning*, pp. 7748–7759, PMLR, 2021.

[226] Y. Liu, Z. Xu, G. Wang, K. Chen, B. Li, X. Tan, J. Li, L. He, and S. Zhao, "Delightfultts: The microsoft speech synthesis system for blizzard challenge 2021," *arXiv preprint arXiv:2110.12612*, 2021.

[227] S. Liu, D. Su, and D. Yu, "Diffgan-tts: High-fidelity and efficient text-to-speech with denoising diffusion gans," *arXiv preprint arXiv:2201.11972*, 2022.

[228] A. Abbas, T. Merritt, A. Moinet, S. Karlapati, E. Muszynska, S. Slangen, E. Gatti, and T. Drugman, "Expressive, variable, and controllable duration modelling in tts," *arXiv preprint arXiv:2206.14165*, 2022.

[229] G. Liu, Y. Zhang, Y. Lei, Y. Chen, R. Wang, Z. Li, and L. Xie, "Promptstyle: Controllable style transfer for text-to-speech with natural language descriptions," *arXiv preprint arXiv:2305.19522*, 2023.

[230] Z. Jiang, J. Liu, Y. Ren, J. He, Z. Ye, S. Ji, Q. Yang, C. Zhang, P. Wei, C. Wang, *et al.*, "Mega-tts 2: Boosting prompting mechanisms for zero-shot speech synthesis," in *The Twelfth International Conference on Learning Representations*, 2024.

[231] Y. Lee, I. Yeon, J. Nam, and J. S. Chung, "Voiceldm: Text-to-speech with environmental context," in *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 12566–12571, 2024.

[232] Y. Gu, Y. Bian, G. Lei, C. Weng, and D. Su, "Durian-e: Duration informed attention network for expressive text-to-speech synthesis," *arXiv preprint arXiv:2309.12792*, 2023.

[233] R. Shimizu, R. Yamamoto, M. Kawamura, Y. Shirahata, H. Doi, T. Komatsu, and K. Tachibana, "Promptttts++: Controlling speaker identity in prompt-based text-to-speech using natural language descriptions," in *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 12672–12676, IEEE, 2024.

[234] M. Kim, S.-W. Chung, Y. Ji, H.-G. Kang, and M.-S. Choi, "Speak in the scene: Diffusion-based acoustic scene transfer toward immersive speech generation," *arXiv preprint arXiv:2406.12688*, 2024.

[235] S. Ji, Z. Jiang, H. Wang, J. Zuo, and Z. Zhao, "Mobilespeech: A fast and high-fidelity framework for mobile zero-shot text-to-speech," *arXiv preprint arXiv:2402.09378*, 2024.

[236] Z. Wang, Y. Wang, M. Li, and H. Huang, "Artspeech: Adaptive text-to-speech synthesis with articulatory representations," in *Proceedings of the 32nd ACM International Conference on Multimedia*, pp. 535–544, 2024.

[237] Y. Xiao, X. Wang, X. Tan, L. He, X. Zhu, S. Zhao, and T. Lee, "Contrastive context-speech pretraining for expressive text-to-speech synthesis," in *Proceedings of the 32nd ACM International Conference on Multimedia*, pp. 2099–2107, 2024.

[238] J. Yeom, H. Kim, J. Choi, C. H. Lee, N. Park, and S. Yoon, "Voiceguider: Enhancing out-of-domain performance in parameter-efficient speaker-adaptive text-to-speech via autoguidance," *arXiv preprint arXiv:2409.15759*, 2024.

[239] Y. A. Li, X. Jiang, C. Han, and N. Mesgarani, "Styletts-zs: Efficient high-quality zero-shot text-to-speech synthesis with distilled time-varying style diffusion," *arXiv preprint arXiv:2409.10058*, 2024.

[240] R. Yamamoto, Y. Shirahata, M. Kawamura, and K. Tachibana, "Description-based controllable text-to-speech with cross-lingual voice control," *arXiv preprint arXiv:2409.17452*, 2024.

[241] N. Park, H. Kim, C. H. Lee, J. Choi, J. Yeom, and S. Yoon, "Nanovoice: Efficient speaker-adaptive text-to-speech for multiple speakers," *arXiv preprint arXiv:2409.15760*, 2024.

[242] S. He, R. Liu, and H. Li, "Multi-source spatial knowledge understanding for immersive visual text-to-speech," *arXiv preprint arXiv:2410.14101*, 2024.

[243] D. Stanton, Y. Wang, and R. Skerry-Ryan, "Predicting expressive speaking style from text in end-to-end speech synthesis," in *2018 IEEE Spoken Language Technology Workshop (SLT)*, pp. 595–602, IEEE, 2018.

[244] C. Yu, H. Lu, N. Hu, M. Yu, C. Weng, K. Xu, P. Liu, D. Tuo, S. Kang, G. Lei, *et al.*, "Durian: Duration informed attention network for multimodal synthesis," *arXiv preprint arXiv:1909.01700*, 2019.

[245] R. Valle, K. Shih, R. Prenger, and B. Catanzaro, "Flowtron: an autoregressive flow-based generative network for text-to-speech synthesis," *arXiv preprint arXiv:2005.05957*, 2020.

[246] J. Betker, "Better speech synthesis through scaling," *arXiv preprint arXiv:2305.07243*, 2023.

[247] S. Ji, J. Zuo, M. Fang, Z. Jiang, F. Chen, X. Duan, B. Huai, and Z. Zhao, "Textrolspeech: A text style control speech corpus with codec language text-to-speech models," in *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 10301–10305, IEEE, 2024.

[248] D. Yang, J. Tian, X. Tan, R. Huang, S. Liu, X. Chang, J. Shi, S. Zhao, J. Bian, X. Wu, *et al.*, "Uniaudio: An audio foundation model toward universal audio generation," *arXiv preprint arXiv:2310.00704*, 2023.

[249] J. Kim, K. Lee, S. Chung, and J. Cho, "Clam-tts: Improving neural codec language model for zero-shot text-to-speech," *arXiv preprint arXiv:2404.02781*, 2024.

[250] E. Casanova, K. Davis, E. Gölge, G. Göknar, I. Gulea, L. Hart, A. Aljafari, J. Meyer, R. Morais, S. Olayemi, *et al.*, "Xtts: a massively multilingual zero-shot text-to-speech model," *arXiv preprint arXiv:2406.04904*, 2024.

[251] Y. Zhang, G. Liu, Y. Lei, Y. Chen, H. Yin, L. Xie, and Z. Li, "Promptspeaker: Speaker generation based on text descriptions," in *2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pp. 1–7, IEEE, 2023.

[252] D. Lyth and S. King, "Natural language guidance of high-fidelity text-to-speech with synthetic annotations," *arXiv preprint arXiv:2402.01912*, 2024.

[253] H. Liu, Z. Chen, Y. Yuan, X. Mei, X. Liu, D. Mandic, W. Wang, and M. D. Plumbley, "Audioldm: Text-to-audio generation with latent diffusion models," *arXiv preprint arXiv:2301.12503*, 2023.

[254] B. Byrne, K. Krishnamoorthi, C. Sankar, A. Neelakantan, D. Duckworth, S. Yavuz, B. Goodrich, A. Dubey, A. Cedilnik, and K.-Y. Kim, "Taskmaster-1: Toward a realistic and diverse dialog dataset," *arXiv preprint arXiv:1909.05358*, 2019.

[255] J. Kahn, M. Riviere, W. Zheng, E. Kharitonov, Q. Xu, P.-E. Mazaré, J. Karadayi, V. Liptchinsky, R. Collobert, C. Fuegen, *et al.*, "Librilight: A benchmark for asr with limited or no supervision," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 7669–7673, 2020.

[256] Y. Shi, H. Bu, X. Xu, S. Zhang, and M. Li, "Aishell-3: A multi-speaker mandarin tts corpus and the baselines," *arXiv preprint arXiv:2010.11567*, 2020.

[257] K. Zhou, B. Sisman, R. Liu, and H. Li, "Emotional voice conversion: Theory, databases and esd," *Speech Communication*, vol. 137, pp. 1–18, 2022.

[258] G. Chen, S. Chai, G. Wang, J. Du, W.-Q. Zhang, C. Weng, D. Su, D. Povey, J. Trmal, J. Zhang, *et al.*, "Gigaspeech: An evolving, multi-domain asr corpus with 10,000 hours of transcribed audio," *arXiv preprint arXiv:2106.06909*, 2021.

[259] B. Zhang, H. Lv, P. Guo, Q. Shao, C. Yang, L. Xie, X. Xu, H. Bu, X. Chen, C. Zeng, *et al.*, "Wenetspeech: A 10000+ hours multi-domain mandarin corpus for speech recognition," in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 6182–6186, 2022.

[260] K. Lee, K. Park, and D. Kim, "Dailytalk: Spoken dialogue dataset for conversational text-to-speech," in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1–5, 2023.

[261] Q. Yang, J. Zuo, Z. Su, Z. Jiang, M. Li, Z. Zhao, F. Chen, Z. Wang, and B. Huai, "Mscenespeech: A multi-scene speech dataset for expressive speech synthesis," *arXiv preprint arXiv:2407.14006*, 2024.

[262] Z. Jin, J. Jia, Q. Wang, K. Li, S. Zhou, S. Zhou, X. Qin, and Z. Wu, "Speechcraft: A fine-grained expressive speech dataset with natural language description," in *Proceedings of the 32nd ACM International Conference on Multimedia*, pp. 1255–1264, 2024.

[263] J. Kominek, T. Schultz, and A. W. Black, "Synthesizer voice quality of new languages calibrated with mean mel cepstral distortion.," in *SLTU*, pp. 63–68, 2008.

[264] A. Rix, J. Beerends, M. Hollier, and A. Hekstra, "Perceptual evaluation of speech quality (pesq)-a new method for speech quality assessment of telephone networks and codecs," in *2001 IEEE International Conference on Acoustics, Speech, and Signal Processing. Proceedings (Cat. No.01CH37221)*, vol. 2, pp. 749–752, 2001.

[265] Wikipedia contributors, "Word error rate — Wikipedia, the free encyclopedia." https://en.wikipedia.org/w/index.php?title=Word_error_rate&oldid=1260041897, 2024. [Online; accessed 7-December-2024].

[266] Wikipedia contributors, "Mean opinion score — Wikipedia, the free encyclopedia." [Online; accessed 7-December-2024].

[267] P. C. Loizou, "Speech quality assessment," in *Multimedia analysis, processing and communications*, pp. 623–654, Springer, 2011.

[268] J. Barker, S. Watanabe, E. Vincent, and J. Trmal, "The fifth'chime'speech separation and recognition challenge: dataset, task and baselines," *arXiv preprint arXiv:1803.10609*, 2018.

[269] X. An, F. K. Soong, and L. Xie, "Disentangling style and speaker attributes for tts style transfer," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 30, pp. 646–658, 2022.

[270] W. Wang, Y. Song, and S. Jha, "Generalizable zero-shot speaker adaptive speech synthesis with disentangled representations," *arXiv preprint arXiv:2308.13007*, 2023.

[271] X. An, F. K. Soong, S. Yang, and L. Xie, "Effective and direct control of neural tts prosody by removing interactions between different attributes," *Neural Networks*, vol. 143, pp. 250–260, 2021.

[272] H. Zen, V. Dang, R. Clark, Y. Zhang, R. J. Weiss, Y. Jia, Z. Chen, and Y. Wu, "Libritts: A corpus derived from librispeech for text-to-speech," *arXiv preprint arXiv:1904.02882*, 2019.

[273] H. Cho, W. Jung, J. Lee, and S. H. Woo, "Sane-tts: stable and natural end-to-end multilingual text-to-speech," *arXiv preprint arXiv:2206.12132*, 2022.

[274] A. Magueresse, V. Carles, and E. Heetderks, "Low-resource languages: A review of past work and future challenges," *arXiv preprint arXiv:2006.07264*, 2020.

[275] A. Farahani, S. Voghoei, K. Rasheed, and H. R. Arabnia, "A brief review of domain adaptation," *Advances in data science and information engineering: proceedings from ICDATA 2020 and IKE 2020*, pp. 877–894, 2021.

[276] Y. Rong and L. Liu, "Seeing your speech style: A novel zero-shot identity-disentanglement face-based voice conversion," *arXiv preprint arXiv:2409.00700*, 2024.

[277] J. Lu, B. Sisman, R. Liu, M. Zhang, and H. Li, "Visualtts: Tts with accurate lip-speech synchronization for automatic voice over," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 8032–8036, 2022.

[278] T. Zhang, "Deepfake generation and detection, a survey," *Multimedia Tools and Applications*, vol. 81, no. 5, pp. 6259–6276, 2022.

[279] Y. Nirkin, Y. Keller, and T. Hassner, "Fsgan: Subject agnostic face swapping and reenactment," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 7184–7193, 2019.

[280] S. Arik, J. Chen, K. Peng, W. Ping, and Y. Zhou, "Neural voice cloning with a few samples," *Advances in Neural Information Processing Systems*, vol. 31, 2018.