

CampNet: Context-Aware Mask Prediction for End-to-End Text-Based Speech Editing

Tao Wang, Jiangyan Yi, *Member, IEEE*, Ruibo Fu, *Member, IEEE*, Jianhua Tao, *Senior Member, IEEE*, and Zhengqi Wen

Abstract—The text-based speech editor allows the editing of speech through intuitive cutting, copying, and pasting operations to speed up the process of editing speech. However, the major drawback of current systems is that edited speech often sounds unnatural due to cut-copy-paste operation. In addition, it is not obvious how to synthesize records according to a new word not appearing in the transcript, which often needs the help of text-to-speech (TTS) and voice conversion (VC) technology at the same time. This paper first proposes a novel end-to-end text-based speech editing method called context-aware mask prediction network (CampNet). The model can simulate the text-based speech editing process by randomly masking part of speech and then predicting the masked region by sensing the speech context. It can solve unnatural prosody in the edited region and synthesize the speech corresponding to the unseen words in the transcript. Secondly, for the possible operation of text-based speech editing, we design three text-based operations based on CampNet: deletion, insertion, and replacement. These operations can cover various situations of speech editing. Thirdly, to synthesize the speech corresponding to long text in insertion and replacement operations, a word-level autoregressive generation method is proposed, which can synthesize the speech of arbitrary length text. Fourthly, we propose a speaker adaptation method using only one sentence for CampNet and explore the ability of few-shot learning based on CampNet, which provides a new idea for speech forgery tasks. The subjective and objective experiments¹ on VCTK and LibriTTS datasets show that the speech editing results based on CampNet are better than TTS technology, manual editing, and VoCo method (the combination of TTS and VC). We also conduct detailed ablation experiments to explore the effect of the CampNet structure on its performance. Finally, the experiment shows that speaker adaptation with only one sentence can further improve the naturalness of speech editing for one-shot learning.

Index Terms—text-based speech editing, text-to-speech, mask prediction, coarse-to-fine decoding, one-shot learning

I. INTRODUCTION

The rapid development of internet has accelerated the transmission of information. There are various media for us to learn, entertain and communicate: movies, podcasts, YouTube videos, interactive online education, etc. The production of these media is often inseparable from speech editing.

T. Wang is with the National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Science, Beijing 100190, China, and also with the School of Artificial Intelligence, University of Chinese Academy of Sciences, Beijing 100190, China (e-mail: tao.wang@nlpr.ia.ac.cn).

J. Yi, R. Fu, J. Tao and Z. Wen are with the National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences, Beijing 100190, China (e-mail: {jiangyan.yi, ruibo.fu, jhtao, zqwen}@nlpr.ia.ac.cn).

Corresponding Author: Jiangyan Yi, Ruibo Fu, Jianhua Tao. E-mail: {jiangyan.yi, ruibo.fu, jhtao}@nlpr.ia.ac.cn.

¹Examples of generated speech can be found at <https://hairuo55.github.io/CampNet>.

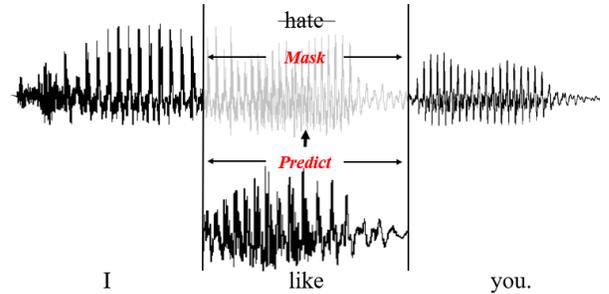


Fig. 1. The replacement operation of text-based speech editing. The operation can be divided into two steps: first, masking the region to be edited and then predicting new speech according to the modified text and speech context. Deletion and insertion operations in text-based speech editing can also be divided into masking and prediction processes, as described in Sec. II-C.

Typical speech editing interfaces [1] present a visualization of the speech such as waveform and/or spectrogram and provide the user with standard select, cut, copy, paste, and volume adjustment, which are applied to the waveform itself. Some advanced operations such as time-stretching, pitch bending, and de-noising are also supported. Such tools provide a great convenience for media producers and have a wide range of application scenarios [2].

Some state-of-the-art systems allow the editor to perform select, cut, and paste operations in the text transcript of the speech and apply the changes to the waveform accordingly, which is called text-based speech editing [3]. In many cases, it is useful to delete some words, insert a new word or phrase in the speech editing process, such as deleting dirty words, replacing a misspoken word or inserting an adjective as emphasis. Text based speech editing can directly modify the speech waveform by deleting, inserting or replacing the target word in the text. However, it mainly faces two problems. One is that the edited speech often sounds unnatural because the edited region does not match the prosody of speech context. (e.g., mismatches in intonation, stress, or rhythm) [4]. Another is that the interfaces do not support the ability to synthesize new words not appearing in the transcript [3]. There are a series of studies on these problems.

For the first problem of prosody mismatch, the main reason is that the prosody of edited speech does not match the context of the original speech. Due to the continuity of human speech, there will be different prosody in different scenes. Directly concatenating the speech in different scenes will lead to problems such as discontinuity of fundamental frequency (F0) and background mismatch. To solve this problem, speech manipulation includes fast pitch-shifting and time-stretching

techniques is applied, which include Time-Domain Pitch-Synchronous Overlap-and Add (TD-PSOLA) [5], WORLD [6], and STRAIGHT [7]. These methods are efficient and suitable for real-time interactive applications but will produce audible artifacts [8, 9]. Neural vocoders such as WaveNet [10], etc. [11–16] can obtain higher perceptual quality than traditional methods, but can not perform context-aware generation for text-based speech editing. To generate prosodies that sound natural in context with any preceding or following speech, context-aware prosody correction [4] is applied to modify the prosodic information of the target segment. In this method, the prosodic information is predicted by neural network, then prosodic modification is realized by applying the TD-PSOLA algorithm [5], followed by de-noising and de-reverberation [17]. This method combines neural network and digital signal processing method, which can effectively improve the speech quality. However, an obvious limitation of this system is that the words to insert or replace may not be found in the available speech data of the target speaker, which limits the application in the field of text-based speech editing.

The second problem is that new words that do not appear in transcripts could not be synthesized. It is easy for a person to type a new word not appearing in the transcript, but it is not obvious how to synthesize the corresponding speech. Of course, it is possible to record new audio with missing words, but it needs to access the original voice talent [3], which will bring great difficulties to the speech editing process. With the rapid development of deep learning in the task of speech generation, the speech synthesized by machines can be comparable to humans, such as the works Tacotron [18, 19] and WaveNet [10] in the field of TTS. Besides, some transfer learning works based on TTS can generate the speech of the target speaker, such as global style token (GST) [20], etc [21–23]. However, these methods are sentence level generation, and it is impossible to edit the specific words in the synthesis speech. To achieve text-based speech editing, the previous work was completed with the help of TTS and voice conversion (VC) system [24, 25], which is called VoCo [3, 26]. The key idea of VoCo is to synthesize the inserted word using a similar TTS voice (e.g., having correct gender) and then modify it to match the target speaker using the VC model [27](making an utterance by one person sound as if it was made by another). Since each module is independent of the other, it will accumulate errors and bring difficulties to the construction of the system. EditSpeech [28] is developed upon a neural TTS framework. It uses force alignment technology to align text with speech first, and then uses partial inference and bidirectional fusion to incorporate the contextual information related to the edited region. This method does not need pipeline structure, but introduces a priori alignment information to realize the mapping between edited text and target region, so as to avoid the unnatural phenomenon caused by cut-copy-paste operation.

Based on the analysis of the above two problems, it can be found that the operation of cut-copy-paste leads to the unnatural phenomenon of speech editing results. There are two reasons. First, it isn't easy to synthesize the edited region in combination with the speech context. Second, during the

process of deleting, cutting or pasting some voice clips, it is easy to occur an unnatural connection. Different from the cut-copy-paste operation, we view text-based speech editing as two steps. We take the replacement operation as an example, as shown in Fig. 1. The replacement process can be divided into two steps. First, we mask part of the original speech that needs to be edited. Then the masked region is predicted according to the modified transcription and speech context. In fact, some other operations of text-based speech editing, such as deletion and insertion, can also be viewed as the process of masking and prediction, which will be described in detail in Sec. II-C. The advantage of this view is that text-based speech editing can be described by an end-to-end model, which can directly generate edited speech according to the context and ensure the natural prosody of speech.

This paper describes approaches to text-based speech editing that can automatically delete, replace and insert the speech at word level by editing the transcription. The approaches can avoid the unnatural problems caused by cutting and pasting in traditional methods and synthesize speech matching with the context in an end-to-end form. Firstly, the context-aware mask prediction network (CampNet) is proposed to simulate the process of text-based speech editing. Secondly, three text-based speech editing operations based on CampNet are designed: deletion, replacement, and insertion. Thirdly, a word-level autoregressive generation method is proposed to improve the editing length. Fourthly, a transfer learning method using only one sentence is proposed, which can further improve the naturalness of the model and provide a new idea for speech forgery.

Overall, the main contributions of this paper are:

- We propose CampNet for the text-based speech editing task, which avoids the unnatural phenomenon caused by “cut-copy-paste” operation in the traditional method and can synthesize new words not appearing in the transcript. To our best knowledge, CampNet is the first text-based speech editing model that can be trained in end-to-end form without duration information. The model uses mask prediction to simulate the text-based speech editing process (Sec. II-A), and uses coarse-to-fine decoder to improve the perception of speech context (Sec. II-B).
- Based on CampNet, we design three speech editing operations, corresponding to delete, replace and insert operations, respectively (Sec. II-C). These operations can comprehensively cover different kinds of situations that text-based speech editing can face.
- In the text-based insertion and replacement operations, to synthesize the speech corresponding to long text, a word-level autoregressive generation method is proposed (Sec. II-D), which can synthesize the speech of arbitrary length text by using CampNet.
- We propose a one-sentence speaker adaptation method for the CampNet, further boosting the performance for one-shot learning, and providing a new idea for speech forgery. Experiments show that this method can obtain better speech similarity for unseen speakers with only one sentence than TTS and VC systems (Sec. II-E).

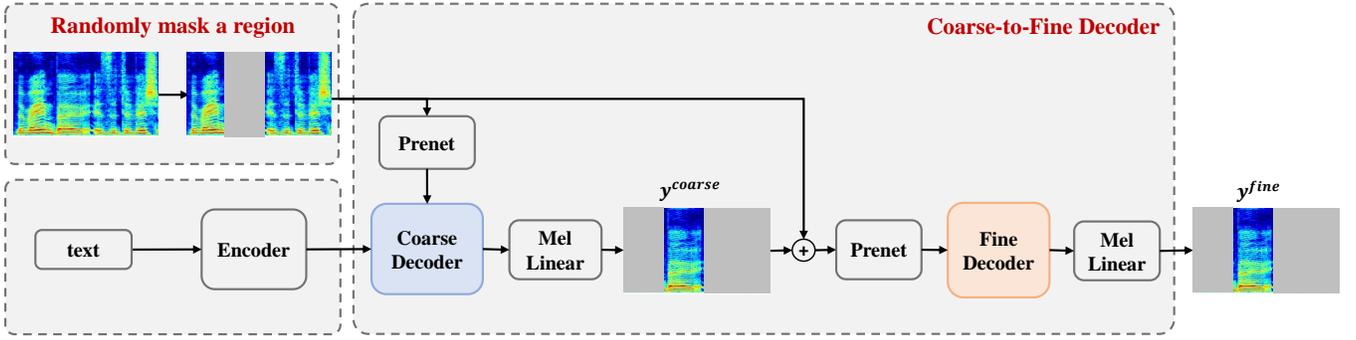


Fig. 2. Structure of Context-Aware Mask Prediction Network (CampNet). To simulate text-based speech editing, the model randomly masks a region of speech at the training stage and then predicts the masked speech according to the corresponding text and speech context.

This paper is structured as follows. The proposed CampNet, its operations for text-based speech editing, word-level autoregressive generation method, and its transfer learning method are described in Section II. After explaining the experiments and results in Section III and Section IV, we draw a conclusion in Section V.

II. CONTEXT-AWARE MASK PREDICTION NETWORK

The context-aware mask prediction network (CampNet) consists of two processing stages, as shown in Fig. 2: encoder and decoder. First, the encoder module processes the input sentence and converts it into a hidden representation. This representation is used to guide the decoder to predict the acoustic feature of the edited speech. Second, a random region of acoustic features is masked as the ground truth to condition the decoder at the decoding stage. The decoder is divided into two steps. The first step is to learn the alignment information between the masked ground truth and the text representation through the multi-head attention mechanism and predict coarse acoustic features. Then, the second step is to predict finer acoustic features based on the coarse acoustic features and original speech context, which can further fuse the context information of speech to make the predicted speech more natural. We call the process of masking part of the acoustic features and predicting the mask region as the "context-aware mask prediction".

In the section, we will first introduce CampNet. Secondly, we will show how to use CampNet for text-based speech editing tasks, including delete, replace and insert operations. Thirdly, the word-level autoregressive generation method is proposed. Finally, we will present the few-shot and one-shot learning methods based on CampNet.

A. Context-Aware Mask Prediction

The task of end-to-end text-based speech editing model is to modify part of the original speech to match the edited transcription. Given the source acoustic features $y = (y_1, \dots, y_n, \dots, y_m, \dots, y_T)$ and its transcription sequences $x = (x_1, \dots, x_a, \dots, x_b, \dots, x_M)$, where (x_a, \dots, x_b) is aligned with acoustic features (y_n, \dots, y_m) . When (x_a, \dots, x_b) in transcription x is edited and the new transcription x' is $(x_1, \dots, x'_a, \dots, x'_b, \dots, x_M)$, the target acoustic feature is $y' = (y_1, \dots, y'_n, \dots, y'_m, \dots, y_T)$. We

assume that the length of the new speech in the editing region is consistent with that of the original region in the training stage. If it is inconsistent at the inference stage, an additional duration model [29] can be used to predict the duration of edited words. Then we can add or delete some fragments on the original part to achieve consistent length. The purpose of this assumption is to ensure that there is no mismatch between the model in the training stage and the test stage. So, the problem of text-based speech editing can be formulated in terms of estimating the condition probability $P(y'|y, x, x'; \theta)$ of the target acoustic feature, and θ is corresponding model parameters, where

$$P(y'|y, x, x'; \theta) = P(y'_n, \dots, y'_m | y, x, x'; \theta) \quad (1)$$

Since (y'_n, \dots, y'_m) is independent of text sequences x and related to text x' , we can get the formula from Eq. 1:

$$P(y'|y, x, x'; \theta) = P(y'_n, \dots, y'_m | y, x'; \theta) \quad (2)$$

In addition, it can be observed that (y'_n, \dots, y'_m) and (y_n, \dots, y_m) have the same position, and their contents are different. To remain the position information of edited region (y'_n, \dots, y'_m) and not be interfered by the content information of (y_n, \dots, y_m) , (y_n, \dots, y_m) can be replaced with a new token $\langle \text{mask} \rangle$, and the masked source acoustic features y_{mask} can be expressed as:

$$y_{\text{mask}} = (y_1, \dots, \langle \text{mask} \rangle, \dots, \langle \text{mask} \rangle, \dots, y_T) \quad (3)$$

Where the n position in y_{mask} is the starting point of the mask, and m is the ending point of the mask.

Then the condition probability $P(y'|y, x, x'; \theta)$ can be formulated as:

$$P(y'|y, x, x'; \theta) = P(y'_n, \dots, y'_m | y_{\text{mask}}, x'; \theta) \quad (4)$$

From the Eq. 4, the task of text-based speech editing can be decomposed into the following two process. First, mask the region of the original speech y that needs to be edited, and get the masked acoustic features y_{mask} . Then, combined with the masked acoustic feature y_{mask} and the edited text sequence x' , neural network is used to predict the edited region (y'_n, \dots, y'_m) . Because x' and y_{mask} have different lengths, and it has been proved that transformer can effectively fuse the context information of sequences with different lengths

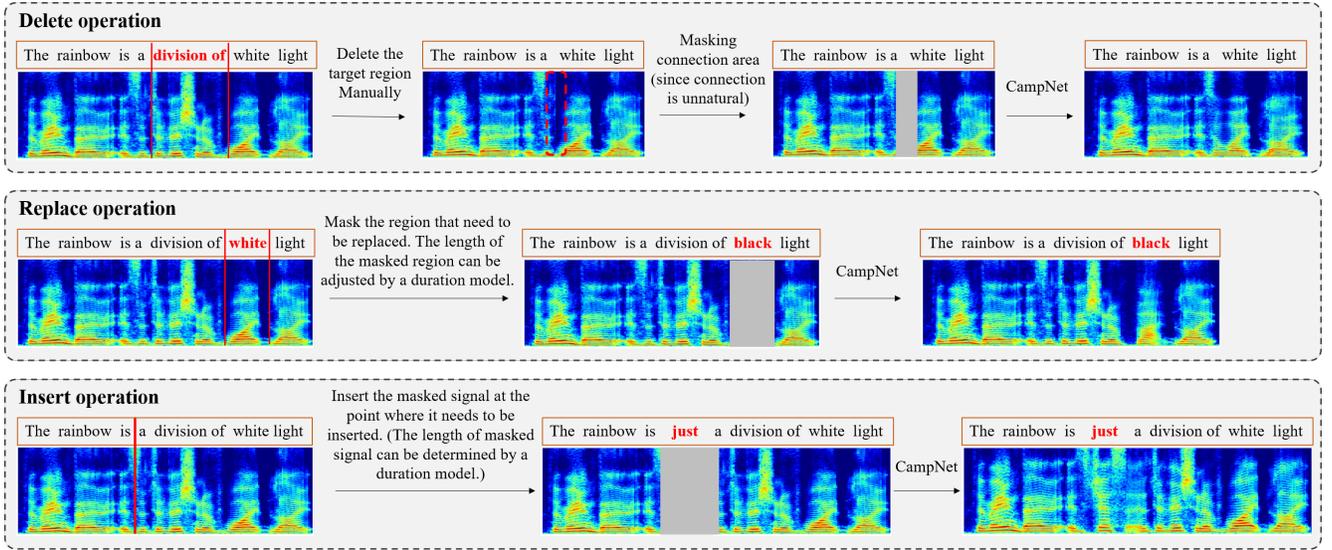


Fig. 3. At the inference stage, three different operations of text-based speech editing based on CampNet are proposed.

[30], an encoder-decoder framework based on transformer is adopt as the structure of CampNet, which is shown in Fig. 2.

In the CampNet, an encoder processes the input sentence $x' = (x_1, \dots, x'_a, \dots, x'_b, \dots, x_M)$ and converts it into a hidden representation in the following way:

$$m = (m_1, m_2, \dots, m_M) = \text{encoder}_{\theta_e}(x') \quad (5)$$

where θ_e denotes the parameters of encoder network.

Since there are only text and speech pairs in the training stage, we randomly mask part of the speech and then use the network to predict the masked part to simulate the text-based speech editing process. The advantage of the masking mechanism is that the training data of speech editing can be well simulated without parallel corpus of speech editing. In addition, due to randomly mask a region during training, a wealth of augmented data can be obtained, improving the robustness of the model. The masked acoustic features y_{mask} are first consumed by a neural network composed of two fully connected layers with ReLU activation [31], named *prenet*. It is responsible for projecting y_{mask} into the same subspace as phoneme embeddings, which can be expressed as:

$$h = (h_1, h_2, \dots, h_T) = \text{prenet}_{\theta_p}(y_{\text{mask}}) \quad (6)$$

where θ_p denotes the parameters of *prenet*.

Finally, combined with m and h hidden features, the final signal is predicted by the decoder, which can be expressed as:

$$(y_{n'}, \dots, y_{m'}) = \text{decoder}_{\theta_d}(h, m) \quad (7)$$

where θ_d denotes the parameters of decoder network.

At the decoding stage, how to perceive the context of speech is the key to synthesize natural speech. We propose a coarse-to-fine decoding method to achieve context-aware, and the details are introduced in the following subsection.

B. Coarse-to-Fine Decoding

To better perceive context information of speech and make the predicted speech more natural, we propose a two-stage

decoder, which is named as “coarse-to-fine decoding”. The structure is a two-stage transformer in series, as shown in Fig. 4. Since our framework is non-autoregressive, when we perform the first decoding (coarse decoding), the model needs to predict all frames of the mask area at one time, and can not fuse the information of the previous frame in the form of autoregression to predict the information of the next frame. Therefore, to better integrate the context information, we perform the second decoding (fine decoding) based on the accumulation of the output of the first-level decoder and the original masked speech information. Then the model can obtain better context information, rather than directly predict the speech information from the masked signal.

Since the transformer can output the whole sequence in parallel, the coarse decoding process can be represented by the product of the probabilities of each frame:

$$P(y^{\text{coarse}} | y_{\text{mask}}, x'; \theta) = \prod_{t=1}^T P(y_t | y_{\text{mask}}, x'; \theta) \quad (8)$$

In this process, we only predict the acoustic features of the masked region, while the output target of the region without masking is the mask value $\langle \text{mask} \rangle$, as shown in Fig. 4. It can be seen from Eq. 8 that each frame of y^{coarse} is output in parallel. Compared with the autoregressive models such as Tacotron and WaveNet, the non-autoregressive model can not fuse the information of the previous frame well. This will cause the naturalness of the synthesized speech to be worse than that of the autoregressive model. To solve this problem, we perform the second decoding, which feeds the sum of y^{coarse} and y_{mask} to a fine decoder to predict finer acoustic features y^{fine} . It can be expressed as:

$$P(y^{\text{fine}} | y^{\text{coarse}}, y_{\text{mask}}; \theta) = \prod_{t=1}^T P(y_t | y^{\text{coarse}} + y_{\text{mask}}; \theta) \quad (9)$$

The second decoding can combine the output of the first-level decoder and the context information of original speech,

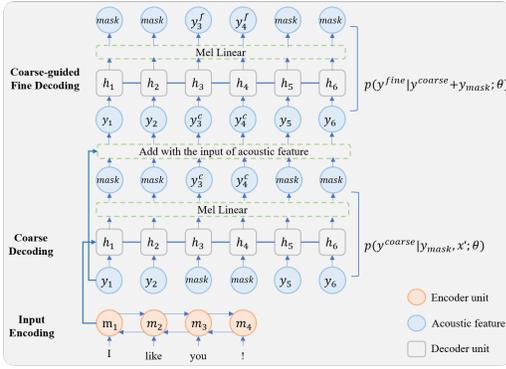


Fig. 4. Structure of coarse-to-fine decoding. We first generate the coarse acoustic features y^{coarse} from text information and masked acoustic features. Then, a fine decoder fills in the missing details based on coarse acoustic features and speech context, and predicts finer acoustic feature y^{fine} .

which is helpful to generate more natural and expressive speech. The target of the second decoder is the same as the first decoder. The advantage of this is that the model can synthesize speech in a non-autoregressive way, and better integrate context information in the form of two-level decoding.

C. Speech Editing Operations Based on CampNet

With a pre-trained CampNet model, some operations of speech editing, such as deletion, insertion, and replacement, can be carried out. The operations are shown in Fig. 3. In this section, we will introduce these in detail.

1) *Delete operation*: The deletion operation allows the user to remove a region of speech waveform that corresponds to certain specified words. We divide the process into three steps. The first step is to manually delete the target region and the corresponding words in the text. Due to manual deletion, unnatural phenomena will appear at the connection, such as fundamental frequency discontinuity. To repair this problem, there are two solutions. One is convenient but empirical, and the other is accurate but requires additional duration model. For the former, we can take the connection point as the center and mask the left and right of speech in a small range. Then, we input the masked speech and the text after deleting the target word into CampNet to re-predict the masked region. The reason for masking a small range at the connection is to re-predict the pronunciation of word-final of the previous word and word-initial of the next word. Since masking a small area at the connection is empirical, another more accurate method is to use an additional duration model to predict the pronunciation duration of word-final of the previous word, word-initial of the next word, as well as the pause information of the two words. Then mask a specific range according to this information. Users can adopt one of two methods according to the needs of different scenarios.

2) *Replace operation*: The replace operation allows the user to replace a fragment of speech with another speech. The operation can be divided into two cases. One is that the length of the replaced segment is close to the target pronunciation. The other is that there is a large gap between them. For the former, it can be divided into two steps. The first step is to

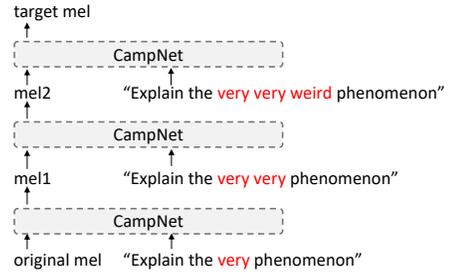


Fig. 5. An example of word-level autoregressive generation method. The figure shows the process of inserting a speech with the text “very very weird” (marked in red) into the original speech.

define the word boundary to be replaced, mask it according to the word boundary and then modify the text. It is worth noting that the range of masking can be larger than the actual boundary when masking. In this way, the model can learn more natural connections. The second step is to input the masked speech and the modified text into CampNet. The model will predict the replaced speech according to the modified text.

If there is a big difference between the length of the replaced speech and the original speech, such as adding some words or deleting some words, a pre-training duration model can be used to predict the length of the replaced region. The duration model is widely used in traditional TTS task[29]. Here, we use the duration model to obtain the speech length of the replaced word. Then according to the predicted length, the masked region can be added or deleted some fragment to ensure the consistency of the duration.

3) *Insert operation*: The insert operation allows the user to insert a speech into the edited speech. This operation is similar to the replacement operation. Firstly, we can use the pre-trained duration model to predict the duration of the words to be inserted. Then insert the masked signal with the predicted length into the original speech. Finally, input the modified text and speech into CampNet, and predict the inserted speech.

It is worth mentioning that when we insert or replace some words, to make the pronunciation of the edited words more natural with adjacent words, we can mask part of the pronunciation of the adjacent words appropriately, such as the word-final of the previous word, and word-initial of the next word. Then, CampNet is used to re-predict the pronunciation of the masked area of these adjacent words, so as to make the prosodic connection more natural.

D. Word-level Autoregression Generation Based on CampNet

In the replacement or insertion operation of text-based speech editing in Sec. II-C, it is easy to face the need to generate speech with long text (For example, the text to be replaced has many words, or the inserted speech has many words). Since we only mask a small part of speech in the training process, when we need to generate speech with long text in the inference stage, the training and inference stages do not match, and the performance of the model will be poor. Therefore, in order to enable CampNet to synthesize long text speech in the inference stage, we propose a word-level generation method, which can effectively solve this problem.

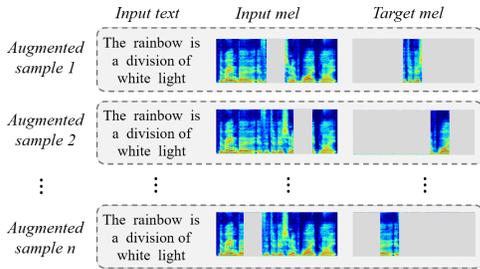


Fig. 6. An example of data augmentation using only one utterance. By randomly masking a region of speech, one sample can be extended to multiple different inputs.

The generation method based on word-level autoregression is shown in the Fig. 5, different from synthesizing the speech corresponding to all texts at one step, we view the generating long text speech as a multi-stage process. First, insert the first word into the original text and synthesize the speech with the modified text. Then, on this basis, insert the second word and cycle in turn until all words are inserted, then the final speech can be obtained. Since this process generates speech word by word, it is called word-level autoregressive generation method.

It should be noted that the word-level autoregressive generation is important for CampNet. This method can make CampNet not limited to the problem of editing length, but can synthesize speech of any length and ensure the stability of generation.

E. Transfer Learning for Text-based Speech Editing

When the target speaker’s corpus is small, synthesizing the target speaker’s voice has always been a research hotspot. Many studies are conducted from the perspective of TTS and VC. Unlike TTS and VC systems, the proposed framework uses the existing speech as input to modify individual words. Compared with synthesizing a whole sentence simultaneously, the task is simpler, so it is easy to get better stability and similarity. In addition, CampNet can obtain the speaker information through the unmasked region without using the speaker embedding, which avoids the error caused by the extraction of the speaker embedding.

After pre-training the CampNet with multi-speaker data, the model can synthesize a new speaker’s voice without transfer learning. We can also fine-tune the parameters for better results, similar to the transfer learning in TTS and VC tasks. However, CampNet can enlarge the training data by randomly masking different positions of speech. Even the training corpus has only one sentence, which can not be achieved by the TTS and VC model. This section will introduce the transfer learning methods based on CampNet for few-shot learning and one-shot learning in detail.

1) *Few-shot learning*: When we need to use a small corpus for speaker adaptation, we only need to transfer the parameters related to speaker features in our proposed framework. Since the encoder only has text input and there is no speaker embedding to control the speech, after pre-training the model with a multi-speaker dataset, we can consider the encoder’s parameters universal. We only need to fine-tune the decoder because the input of the decoder is the acoustic features. After

fine tuning the decoder, the model can predict the acoustic features that better match the target speaker.

2) *One-shot learning*: One-shot learning has always been the difficulty of speech forgery. This section provides a different idea from TTS and VC, which can obtain better similarity. We can directly fine-tune the model with only one sentence to improve the similarity with the target speaker who has never appeared in the training corpus. Because CampNet is based on mask and prediction mechanism, randomly masking acoustic features can effectively augment the training data and improve the robustness of the model, which is shown in Fig. 6. After several steps of fine-tuning the model with one sentence, the model is not easy to overfit, and the performance will be further improved due to the difference of input at each step.

III. EXPERIMENTAL PROCEDURES

A. Dataset and Task

In this section, we conduct experiments on VCTK [32] and LibriTTS [33] corpus to evaluate our proposed method. The VCTK corpus includes speech data uttered by 110 English speakers with different accents. Each speaker reads out about 400 sentences. Specifically, we select four speakers from the VCTK dataset as the test set, and the rest utterances are divided into 90% training set and 10% validation set. We use the training set to train the model. We also randomly select 100 sentences from the LibriTTS corpus as the test set to verify the model’s performance on the cross dataset.

We mainly compare the replacement operation of CampNet and other systems, which is easy to evaluate with the original speech². To ensure that the edited speech of different systems is consistent with the content of the original speech, which facilitates the comparison between objective metrics and subjective metrics, we randomly choose 80 words that span 3 to 10 phonemes from 80 different sentences for each test set. For each sentence, we remove the region of the corresponding words in the speech. Then we use different systems to predict the removed region. All wav files are sampled at 16KHz.

B. Model Details

Acoustic features are extracted with a 10 ms window shift. LPCNet [16] is utilized to extract 32-dimensional acoustic features, including 30-dimensional BFCCs [34], 1-dimensional pitch and 1-dimensional pitch correction parameter. Five systems are compared in our experiments, including TTS technology, the combination of TTS and VC, manual editing, CampNet and actual recording. We use the training set in VCTK to train the LPCNet model.

- **Synth** We train a neural TTS system to synthesize the speech and copy the target region to insert into the edited speech. To make the voice of the synthesized speech as similar as that of the target speaker as possible, Tacotron2 based on global style token (GST [20]) is used as the acoustic model. The structure of Tacotron2 is the same as that in paper [19]. The hidden dimension of the encoder and decoder in Tacotron is 512, and the GST dimension

²Examples of more operations can be found at <https://hairuo55.github.io/CampNet>.

is 128. We use the phoneme as input and output the 32-dims acoustic features extracted by LPCNet. The initial learning rate is set to 1e-3. Adam [35] is used as the optimizer.

- VoCo** The main idea of VoCo is to synthesize the inserted word using a similar TTS voice (e.g., having correct gender) and then modify it to match the target voice using a VC model. To realize the VoCo system, we train the neural TTS system and VC system, respectively. For the TTS system, the configuration is the same as **Synth**. We select two speakers (one male and one female) from the training set as source speakers. Then, the TTS system is used to synthesize speech, and the VC system synthesizes the voice which is similar to the target speaker. We copy the target region from the output speech and insert it into the edited speech. The voice conversion is based on phonetic posteriorgrams (PPGs) [27], which can be applied to non-parallel voice conversion and achieved both high naturalness and high speaker similarity of the converted speech. The 512-dimensional PPGs are extracted from the acoustic model in speaker independent-ASR, which is implemented using the Kaldi toolkit [36] and trained on our 20,000 hours corpus. The voice conversion model's structure follows the structure in paper [27]. Furthermore, to make the voice conversion system have the one-shot ability, we use GST as speaker embedding to train a multi-speaker VC model. The initial learning rate of the VC model is set to 1e-3, and Adam [35] is used as the optimizer. It should be noted that the TTS technology in the original VoCo adopts the unit selection technology, which is complex to realize in VCTK datasets. While, at present, the popular end-to-end TTS technology can also generate natural speech. Therefore, we mainly adopt the idea of realizing speech editing with TTS and VC technology in VoCo, and use Tacotron to realize TTS.
- Edit** We use the editing interface to refine the speech further if it improves on Synth/VoCo.
- Real** The actual recording without modification.
- CampNet** The proposed framework of CampNet is shown in Fig. 2. The structures of encoder, coarse decoder and fine decoder are based on transformer [37], as shown in Fig. 7. We input the phoneme sequence into a 3-layer CNN [38] to learn the context information of the text. Each phoneme has a trainable embedding of 256 dims, and the output of each convolution layer has 256 channels, followed by batch normalization and ReLU activation and a dropout layer as well [31, 39, 40]. The transformer blocks of the encoder and fine decoder are 3. The transformer block of the coarse decoder is 6. The hidden dimension of the transformer is 256. At the training stage, we set the masked region to be 12% of the total speech length. The initial learning rate is 1e-3. Adam [35] is used as the optimizer. All batches are set to 16, and the number of training steps is 2 million.

C. Objective Evaluation Metrics

The quality of edited speech can be evaluated by comparing it with the actual speech using the following metrics. To reflect

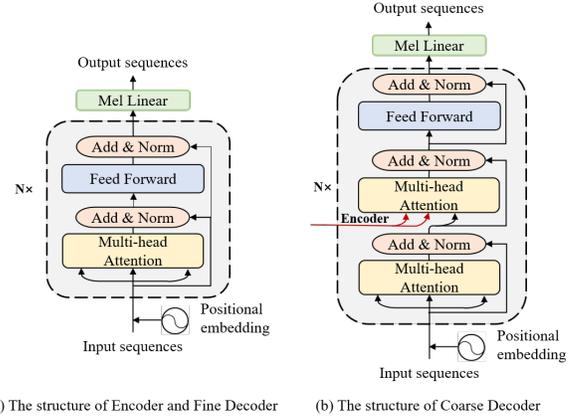


Fig. 7. Structures of encoder, coarse decoder and fine decoder of CampNet.

the overall information of speech (including the information of the editing area and the information of the junction), we calculate the objective metrics of the whole speech. Specifically, since CampNet only synthesizes acoustic features, speech waveforms need to be generated by vocoders. In contrast, other systems paste the target speech directly into the original speech waveform. Therefore, to avoid the influence of vocoder, we also paste the target region predicted by CampNet to the corresponding position of the original speech when calculating the objective metrics. Such objective comparison can be more accurate and fair. In this way, the speech of the unedited area of all systems can be guaranteed to be consistent. In addition, the test set of each system is the same, so it can be ensured that the metrics are only related to the speech of the editing area, and it will not cause unfairness in the comparative experiment.

1) *Mel-Cepstral Distortion (MCD)*: Given two mel-cepstra $\hat{\mathbf{x}} = [\hat{x}_1, \dots, \hat{x}_M]^T$ and $\mathbf{x} = [x_1, \dots, x_M]^T$, we use the mel-cepstral distortion (MCD) [41]:

$$\text{MCD}[\text{dB}] = \frac{10}{\ln 10} \sqrt{2 \sum_{i=1}^M (\hat{x}_i - x_i)^2} \quad (10)$$

to measure their difference. Where M is the order of mel-cepstrum and M is 28 in our implementation. Here, we used the average of the MCDs taken along the DTW [42] path between edited and reference feature sequences as the objective performance measure for each test utterance.

2) *F0-RMSE (Hz)*: For the F0 of speech, following RMSE is applied [43]:

$$\text{F0-RMSE} = 1200 \sqrt{((\log_2(F_r) - \log_2(F_s))^2)} \quad (11)$$

where the subscript r and s represent reference and edited speech, respectively. The F0-RMSE is calculated for each frame, and we use the average of the F0 taken along the DTW path between converted and reference feature sequences.

3) *V/UV error*: The ratio of the number of unmatched U/V frames between reference and edited speech to total frames is calculated as the V/UV error [43]. For two different lengths of speech, we still use the DTW algorithm to align them.

4) *F0-CORR*: We use the correlation coefficient between the edited and reference F0 contours as the objective performance measure to evaluate the F0 of edited speech [43].

TABLE I
OBJECTIVE EVALUATION RESULTS OF SYNTH, VOCO, EDIT AND CAMPNET ON THE TEST SETS OF VCTK AND LIBRITTS

Metrics	VCTK				LibriTTS			
	Synth	VoCo	Edit	CampNet	Synth	VoCo	Edit	CampNet
MCD(dB)	0.594	0.589	0.582	0.380	0.871	0.894	0.870	0.628
F0-RMSE(Hz)	10.463	10.555	10.451	8.637	21.898	2.093	21.308	20.201
V/UV error	1.944	1.843	1.937	1.635	3.916	4.347	3.956	3.675
F0-CORR	0.973	0.972	0.975	0.981	0.940	0.945	0.939	0.954

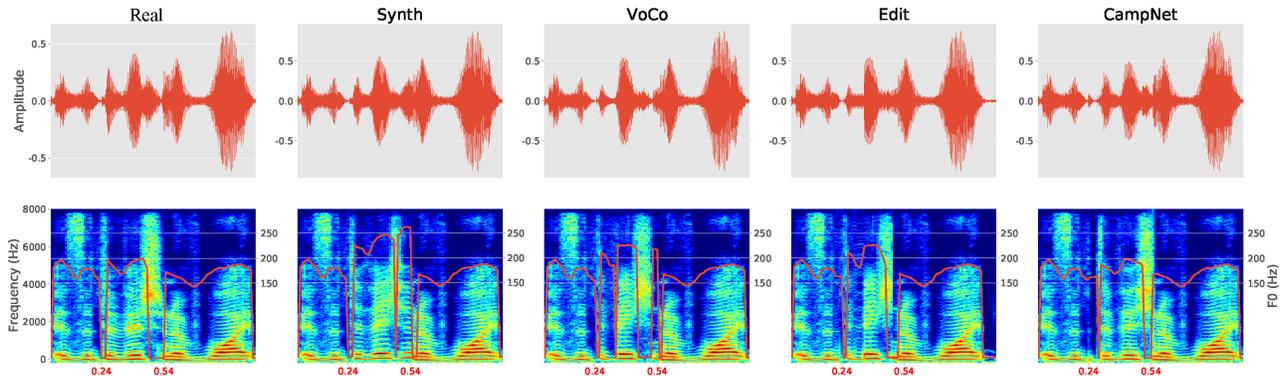


Fig. 8. The waveforms and spectrograms of natural speech and speech edited by different system (the speaker did not appear in the training set). The region marked with time (0.24s ~ 0.54s) is the edited region. The text of the masked region is 'division'.

Since the synthesized and reference speech are not necessarily aligned in time, we computed the correlation coefficient after properly aligning them using the DTW algorithm. If we use $\tilde{\mathbf{y}} = [\tilde{y}_1, \dots, \tilde{y}_{M'}]$ and $\mathbf{y} = [y_1, \dots, y_{M'}]$ to denote the vectors consisting of the elements which is aligned. We can use the correlation coefficient between $\tilde{\mathbf{y}}$ and \mathbf{y}

$$R = \frac{\sum_{m'=1}^{M'} (\tilde{y}_{m'} - \tilde{\varphi})(y_{m'} - \varphi)}{\sqrt{\sum_{m'=1}^{M'} (\tilde{y}_{m'} - \tilde{\varphi})^2} \sqrt{\sum_{m'=1}^{M'} (y_{m'} - \varphi)^2}} \quad (12)$$

where $\tilde{\varphi} = \frac{1}{M'} \sum_{m'=1}^{M'} \tilde{y}_{m'}$ and $\varphi = \frac{1}{M'} \sum_{m'=1}^{M'} y_{m'}$, to measure the similarity between the two F0 contours.

IV. RESULTS

In this section, we first compare the performance of CampNet and some other systems, such as objective metrics, subjective metrics, and operating efficiency. Then, some ablation experiments based on CampNet are conducted. Finally, the ability of few-shot learning and one-shot learning based on CampNet is explored.

A. Comparison between CampNet and Some other Method

This section compares the performance of our proposed **CampNet** with three other speech editing methods, including **Synth**, **VoCo**, **Edit** by objective and subjective evaluations.

First, the objective results on the test sets of VCTK and LibriTTS are listed in Table I. In general, it can be found that the metrics on the two test datasets of CampNet are the best among all the systems. Specifically, in the frequency domain, the CampNet obtained the lowest MCD, which means that human perception would be better. Besides, the F0 has a significant influence on speech perception. We can find that CampNet achieves the best performance in F0-related metrics (F0-RMSE, V/UV error, and F0-CORR). The results show

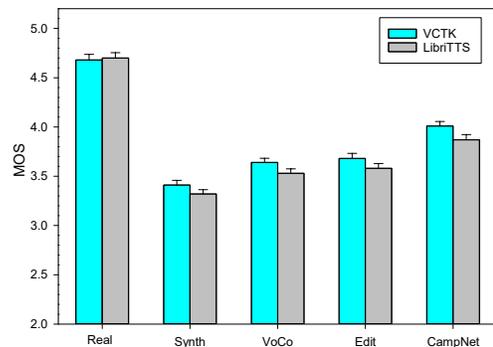


Fig. 9. The MOS score with 95% confidence intervals of the four systems and real speech.

that CampNet can obtain more accurate fundamental frequency information.

Second, we show the waveforms and spectrograms of natural speech and the edited speech generated by different methods. We take the *p225_007.wav* in the VCTK corpus as an example, as shown in Fig. 8. The region marked with time (0.24s ~ 0.54s) is the edited region. It is worth mentioning that the speaker *p225* in the test set does not appear in the training set. It can be found that there will be an unnatural connection in the speech edited by **Synth** and **VoCo** models. There are apparent F0 discontinuities in the frequency domain. In addition, we draw the curve of F0. It can be found that the F0 of speech synthesized by **Synth** and **VoCo** are higher than the original region, while the speech synthesized by CampNet is consistent with the surrounding F0.

Third, subjective evaluations are conducted to compare the performance of CampNet with other systems in terms of the naturalness of edited speech. In this evaluation, twenty utterances in each test set are selected and edited using the proposed method and other systems, including Synth, VoCo, and Edit. It should be noted that all the speakers of the test

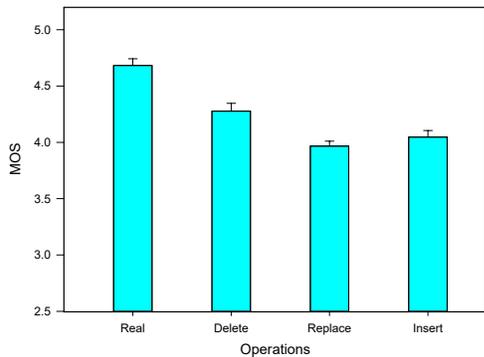


Fig. 10. The MOS score with 95% confidence intervals of different operations and real speech.

data did not appear in the training set. Twenty listeners took part in the evaluation. They are told in advance which word is predicted. The listeners were asked to listen and rate the quality of the restored sentence on a Likert scale [44]: 1 = bad (very annoying), 2 = poor (annoying), 3 = fair (slightly annoying), 4 = good (perceptible but not annoying) and 5 = excellent (imperceptible, almost real). They can play the recording multiple times. Fig. 9 shows the MOS score of each system. The results show that the CampNet is better than the other three systems in each test set. This is also consistent with the previous analysis of objective metrics. In addition, we compare the MOS scores of different operations. Specifically, we have prepared three test sets for the three operations. The speech of each test set has 20 sentences. Twenty listeners took part in the evaluation. They can play the recording multiple times. Fig.10 shows the MOS score of each operation. The results show that the MOS of delete operation is higher. The MOS score of insertion operation and replacement operation is similar, and their gap is small.

In addition, we also compared CampNet with EditSpeech [28] and context-aware prosody correction method [4]. To make the comparison more obvious, we directly use the speech in the demo page of these systems as the standard, and use CampNet to generate the corresponding parallel speech as the comparison. We have put the samples on our demo page (<https://hairuo55.github.io/CampNet>), and we are looking forward to readers' listening.

B. Inference Speed

We evaluate the inference speed of our proposed method with other systems. Since in the Synth and VoCo systems, Tacotron is used as the acoustic model, which significantly impacts the operation efficiency of the whole system. In this section, we compare CampNet with Tacotron and Transformer-TTS. The details of these TTS models are as follows:

- **Tacotron** represents for the TTS model in which the decoder is based on LSTM. The structure details are the same as the acoustic model in the Synth system.
- **Transformer-TTS** represents for TTS model which is based on transformer structure [30]. The structure details are the same as the model in the paper [30]. The number of encoder's blocks is 3, and the number of decoder's blocks is 9, which is consistent with the number of

TABLE II
THE COMPARISON OF INFERENCE SPEED WITH 95% CONFIDENCE INTERVALS FOR CAMPNET, TACOTRON2 AND TRANSFORMER-TTS. THE VALUE OF INFERENCE SPEED INDICATES HOW LONG IT TAKES TO SYNTHESIZE 500 FRAMES OF ACOUSTIC FEATURES.

Model	Params	Inference(s)	Speedup
Tacotron	3.94e7	1.501 ± 0.280	/
CampNet	1.47e7	0.044 ± 0.015	34×
Transformer-TTS	1.52e7	9.460 ± 1.092	/
CampNet	1.47e7	0.044 ± 0.015	215×

decoder blocks of CampNet. The number of hidden layer features is 256.

The evaluation experiments are conducted with 52 Intel Xeon CPU, 512GB memory, and 1 NVIDIA V100 GPU. It is worth mentioning that, during the model's design, we have kept the model parameters as consistent as possible to eliminate their effects. Each model outputs 500 frames of acoustic features for fair comparison. We show the generation speed of acoustic features in Table II. It can be seen that the CampNet speeds up acoustic features generation by 34 times, compared with the Tacotron model. CampNet speeds up acoustic features generation by 215 times, compared with the Transformer-TTS model. It shows that autoregressive generation greatly affects the model's speed. CampNet can synthesize speech in the form of non-autoregressive and effectively improve the synthesis efficiency.

C. The Comparison of Alignment with Tacotron

In this section, we explore the alignment mechanism of the CampNet model, which can help us understand the process of speech editing better. Since the coarse decoder is used to align text and speech, this section explores its attention mechanism.

First, we visualize the alignments of text and speech in the multi-head attention of coarse decoder, and the alignment of local sensitive attention in Tacotron is also visualized for comparison, as shown in Fig. 11. We can find that the alignment of Tacotron is aligned in the whole time step, where each column denotes the attention probabilities corresponding to different encoder states for one decoder step. On the contrary, the alignment of CampNet is only in a small region, that is, the edited region. Its mapping area on the encoder is the edited text. Therefore, CampNet's attention mechanism only focuses on the edited region and automatically finds the text corresponding to the edited region. Compared with the whole alignment in Tacotron, the advantage of this method is that it can make the model pay more attention to the edited region and make the predicted speech consistent with the context. Aligning only a region is simpler than aligning all, which is also in line with our intuitive understanding. Moreover, because text-based speech editing only requires partial alignment, the quality of the training corpus needed is not too high, while the TTS corpus should be as standard as possible.

Fig. 12 shows the alignment results after masking different words in a sentence. It can be seen that the model has good robustness to edit different words in different positions. In addition, the last figure of Fig. 12 shows the alignment of

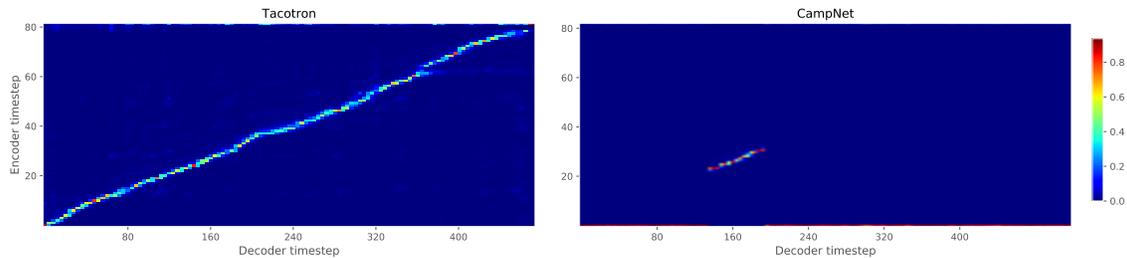


Fig. 11. Comparison of Tacotron alignment and CampNet alignment. The alignment of Tacotron is to align with the text in a complete time step. CampNet only aligns the edited speech region with the edited text.

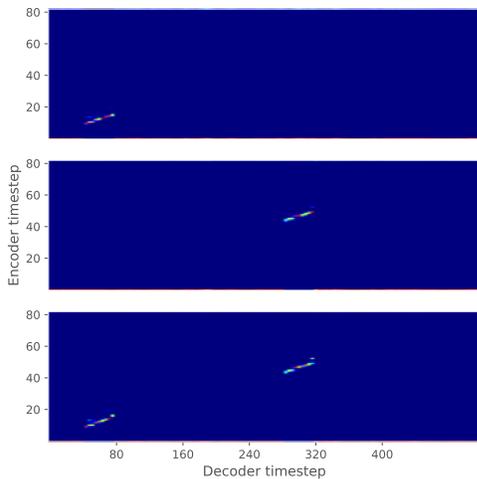


Fig. 12. Alignments after masking the speech segments at different positions. The last one is the alignment of masking the two region at the same time.

two words edited simultaneously during the inference stage. Although we only mask one region of speech at the training stage, we can modify two different regions at the same time at the inference stage. This also verifies that the model can effectively learn the alignment between speech and text and only focus on the edited region.

D. Effectiveness of Coarse-to-Fine Decoding

As introduced in Section II-B, a coarse-to-fine decoding method is proposed to boost performance. To further understand the role of coarse-to-fine decoding, we compare the following three kinds of speech.

- **One-decoder** represents the speech output by the model, which removed the coarse-to-fine decoding of CampNet. We directly use one decoder to output the final speech. To ensure the block number of the decoder are equal to the CampNet, we set the block numbers of the decoder as 9, which is the sum of the coarse decoder and fine decoder.
- **Coarse-decoder** represents the speech $y^{coarser}$ output by the coarse decoder of CampNet.
- **Fine-decoder** represents the speech y^{fine} output by the fine decoder of CampNet.

First, we visualize some spectrograms of original speech and the edited speech after modifying one of the words, as shown in Fig. 13, where 0.24s to 0.88s is the edited region. The speaker of the sample does not appear in the training

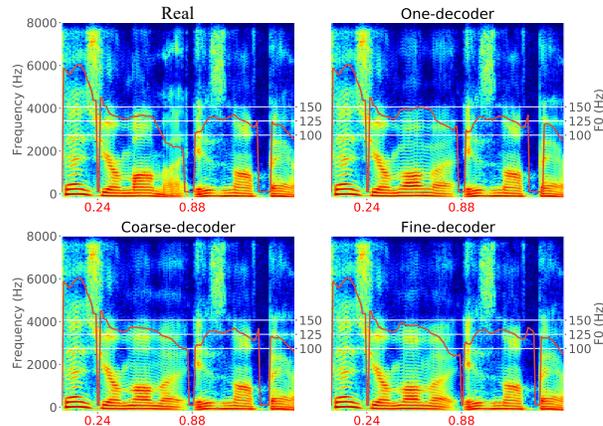


Fig. 13. The spectrograms of the speech generated by different methods. The region marked with time (0.24s ~ 0.88s) is the edited region.

TABLE III
OBJECTIVE EVALUATION RESULTS OF ONE-DECODER, COARSE-DECODER AND FINE-DECODER ON THE TEST SETS OF VCTK AND LIBRITTS.

		Coarser-decoder	One-decoder	Fine-decoder
VCTK	MCD(dB)	0.395	0.388	0.380
	F0-RMSE(Hz)	9.108	8.851	8.637
	V/UV error(%)	1.712	1.686	1.635
	F0-CORR	0.977	0.978	0.981
LibriTTS	MCD(dB)	0.639	0.632	0.628
	F0-RMSE(Hz)	21.650	20.859	20.201
	V/UV error(%)	3.859	3.766	3.675
	F0-CORR	0.945	0.951	0.954

corpus. From the perspective of F0, it can be found that the speech output by the fine decoder is closest to the F0 of the real speech. The F0 of the One-decoder is higher than the original F0. Although the F0 of the Coarse-decoder is different from the real one, after adjusting by the fine decoder, a similar F0 curve is obtained. Besides, we can find that the speech of the fine decoder has a more precise spectrum.

Second, we compare the objective metrics of the three kinds of speech, as shown in Table III. It can be found that the speech of the fine decoder achieves the best performance in all metrics on the two test sets. Specifically, the speech output by the fine decoder is closer to the natural speech in both spectrum and F0, which shows that coarse-to-fine decoding can learn more accurate frequency domain information.

Third, we conduct a subjective ABX test to compare the three kinds of speech. In each subjective test, twenty sentences are randomly selected from the LibriTTS test set. Twenty listeners evaluate each pair of generated speech. The listeners are

TABLE IV
OBJECTIVE EVALUATION RESULTS OF DIFFERENT MASK RATIO AT INFERENCE STAGE ON THE VCTK TEST SET

Metrics	VCTK						LibriTTS					
	M-6%	M-8%	M-10%	M-12%	M-14%	M-16%	M-6%	M-8%	M-10%	M-12%	M-14%	M-16%
MCD(dB)	0.465	0.387	0.391	0.380	0.383	0.398	0.746	0.661	0.631	0.628	0.634	0.650
F0-RMSE(dB)	10.511	9.723	9.407	8.637	9.255	9.114	22.895	22.242	21.049	20.201	21.086	21.820
V/UV error	1.989	1.658	1.610	1.635	1.492	1.750	5.679	4.136	3.635	3.675	4.000	4.259
F0-CORR	0.971	0.976	0.977	0.981	0.978	0.978	0.952	0.943	0.949	0.954	0.948	0.945

TABLE V
AVERAGE PERFORMANCE SCORE(%) ON SPEECH QUALITY AMONG DIFFERENT SYSTEMS, WHERE N/P STANDS FOR "NO PREFERENCE". THE p -VALUES <0.01.

System A	Scores A(%)	N/P Neural(%)	Scores B(%)	System B
CoareDecoder	22.25	47.00	30.75	OneDecoder
FineDecoder	43.25	32.75	24.00	OneDecoder
FineDecoder	44.75	35.00	20.25	CoareDecoder

TABLE VI
OBJECTIVE EVALUATION RESULTS OF DIFFERENT MASK RATIO AT INFERENCE STAGE ON THE VCTK TEST SET

	0.5s	1.0s	1.5s	2.0s	2.5s
MCD(dB)	0.345	0.628	0.960	1.296	1.387
F0-RMSE(dB)	4.577	7.133	10.887	13.403	14.584
V/UV error	1.489	2.783	4.626	6.201	8.175
F0-CORR	0.994	0.986	0.964	0.945	0.931

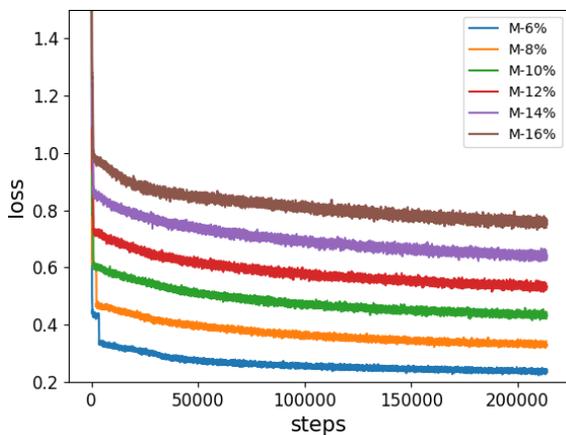


Fig. 14. Comparison of loss functions with different mask ratios at the training stage. The smaller the masked region is, the easier the model is to be optimized and the lower the loss is.

asked to judge which utterance in each pair has better speech quality or no preference in the edited area. They were told in advance which word was predicted. The results are listed in Table V. Fine-decoder outperforms Coarse-decoder and One-decoder, which is consistent with the objective metrics analysis. This result further shows the effectiveness of coarse-to-fine decoding.

It is worth noting that whether the model is improved can be judged through numerical comparison of the objective metrics, but this is not absolutely relevant. More importantly, we should judge it by subjective evaluation. We have put the relevant samples on the demo page (<https://hairuo55.github.io/CampNet/>), and we recommend readers to listen.

E. The Ability of Editing Different Length

As introduced in Section II-E1, we propose the mask prediction method to simulate the speech editing process. The core idea is to randomly mask a speech region at the training stage and then predict the masking region. There may be a mismatch between the length of the masking region at the training stage and the inference stage, which may affect the

model’s performance. This section mainly explores the impact of this problem on the model performance.

First, we explore the influence of different mask ratios during training. Specifically, the mask ratios at the training stage are set to 6%, 8%, 10%, 12%, 14%, and 16%, respectively. The trained models are represented by M-6%, M-8%, M-10%, M-12%, M-14%, and M-16%. All models are trained for 2 million steps with the same structure and hyper-parameters. Fig. 14 compares the change of loss of each model with the increase of training steps. Obviously, with the rise of mask ratios, the loss function becomes larger and larger, which means it is more difficult for the model to predict the masked speech. The smaller the mask ratios, the easier the model is to be optimized. Further, to test the effect of each model, we calculate the objective metrics of each model on the test set, which is shown in Table IV. It can be found that when the mask ratio is set to 12%, the model has the best effect on most indicators (MCD, F0-RMSE, and F0-CORR). Therefore, we use a mask ratio of 12% in the training phase, which performs well in the test set and can also optimize the training loss to a suitable level.

Second, we explore the impact of the edited region’s length at the inference stage in one-step. In order to control the length of the editing area in the test set within the specified length, we directly mask a region in a fixed length and use CampNet to predict the mask region. Referring to Table IV, we set the mask ratio as 12% during the training stage. For all speech in the test set, we take a fixed position of the speech as the starting point, then mask the speech with different lengths, which are 0.5s, 1.0s, 1.5s, 2.0s, and 2.5s, respectively. We input the masked speech into CampNet to predict the masked speech. We calculate the objective metrics of the predicted speech and the real speech, as shown in Table VI. It can be found that with the increase of mask length on the test set, the model’s performance shows a downward trend. The best effect on the test set is the smallest mask length, that is **0.5s**. This indicates that the smaller the mask area is, the closer the synthesized speech is to the real speech. When the mask length exceeds 1.5s, the synthesized speech is quite different from the actual speech in the sense of hearing. Therefore, the length of the editing area should not be too long during the

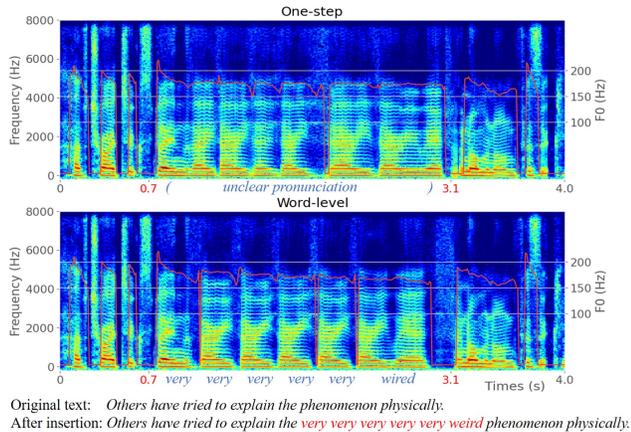


Fig. 15. The spectrograms of the speech generated by different methods. The region marked with red color (0.7s ~ 3.1s) is the insertion area, and its correspondin text is "very very very very very weird". The pronunciation in the insertion area of method One-step is not clear, while the pronunciation of method Word-level is natural and normal.

inference stage, preferably less than 1.5s under the condition that the mask ratio at the training stage is 12%. It is worth mentioning that this length is sufficient to deal with common scenes in text-based speech editing, such as replacing wrong pronunciation or inserting emphasis words.

F. Effectiveness of Word-level Autoregressive Generation Method

As introduced in Section II-D, we propose a word-level autoregressive generation method based on CampNet to face the situation of generating speech with long text. In this section, to explore the ability of word-level autoregressive generation method, we take the insert operation as an example, different from the previous experiments, the length of the words to be replaced is much longer than that used in the previous experiments. The number of words to be inserted is more than 5, and the corresponding speech length is more than 2 seconds. We compared the following two generation methods:

- **One-step** means to generate the speech of all words to be inserted at one step.
- **Word-level** means to generate the speech by using word-level autoregressive generation method, which is introduced in Section II-D.

First, we compare the spectrum generated by the two methods after inserting some words, as shown in Fig. 15. It can be seen that when the word-level autoregressive generation method is adopted, the generated speech is more stable. When multiple words are generated at one step, the pronunciation of the synthesized speech will be unclear. The reason is that only short speech segments is masked in the training stage, which makes it impossible to generate long segments of speech during the inference stage. While, the word-level autoregressive generation method can ensure the matching of mask length in inference and training stage.

Second, since there is no real speech as a comparison when inserting long text, it is not convenient for us to make objective evaluation. Therefore, we conduct a subjective ABX test to

TABLE VII
AVERAGE PREFERENCE SCORE(%) ON SPEECH QUALITY AMONG DIFFERENT SYSTEMS, WHERE N/P STANDS FOR "NO PREFERENCE". THE p -VALUES < 0.01 .

System	Scores	N/P	Scores	System
A	A(%)	Neural(%)	B(%)	B
One-step	5.25	8.00	86.75	Word-level

compare the two methods. In each subjective test, twenty sentences are randomly selected. Twenty listeners evaluate each pair of generated speech. The listeners are asked to judge which utterance in each pair has better speech quality or no preference in the edited area. They were told in advance which word was predicted. The results are listed in Table VII. It is obviously that, in the case of a large number of words to be inserted, the generation method based on word-level autoregression is better than the one-step generation method. For more samples, please refer to our demo page.

G. The Ability of One-shot and Few-shot Learning

As introduced in Section II-E, we propose a transfer learning method based on CampNet. In this section, we explore the ability of CampNet to face few-shot and one-shot learning by comparing the following three models:

- **1-utt_wo_finetune** means to directly edit the speech of an unseen speaker using CampNet without fine-tuning the model.
- **1-utt_w_finetune** means to fine-tune CampNet with one sentence from an unseen speaker before speech editing, which is introduced in Section II-E2.
- **50-utts_w_finetune** means to fine-tune CampNet with 50 sentences from an unseen speaker before speech editing, which is introduced in Section II-E1.

The steps of fine-tuning are five epochs of the dataset used to fine-tune. The fine-tuning datasets are separate from the data used to calculate objective and subjective metrics. We calculate the objective metrics of the synthesized speech and real speech of each system, as shown in Table VIII. We can find that fine-tuning the model with a small amount of corpus can significantly improve the performance in objective metrics. Even if one utterance is used for fine-tuning the model, it can be found that the objective metrics are improved by comparing system **1-utt_w_finetune** and **1-utt_wo_finetune**. This shows the effectiveness of using only one sentence to adapt the model.

Besides, we conduct an ABX test on the three methods. Twenty sentences of each speaker are synthesized by two comparative systems. Twenty listeners evaluate each pair of generated speech. The listeners are asked to judge which utterance in each pair has better speech quality or no preference. The p -value is used to measure the significance of the difference between two systems. The results are listed in Table IX. It can be found that the quality of speech can be further improved through fine-tuning. Furthermore, using only one sentence, the adaptive model is also significantly improved than the non-adaptive model.

TABLE VIII
OBJECTIVE EVALUATION RESULTS OF DIFFERENT METHODS ON THE TEST SETS.

Speaker	Metrics	1-utt_wo_finetune	1-utt_w_finetune	50-utts_w_finetune
P225 (female)	MCD(dB)	0.395	0.273	0.275
	F0-RMSE(Hz)	7.557	4.630	4.614
	V/UV error(%)	1.655	1.325	1.338
	F0-CORR	0.983	0.992	0.992
P226 (male)	MCD(dB)	0.367	0.272	0.271
	F0-RMSE(Hz)	9.406	7.740	6.987
	V/UV error(%)	1.627	1.173	1.137
	F0-CORR	0.958	0.966	0.978

TABLE IX
AVERAGE PREFERENCE SCORES (%) ON SPEECH QUALITY OF DIFFERENT METHODS, WHERE N/P STANDS FOR “NO PREFERENCE”, AND p DENOTES THE p -VALUE OF A t -TEST

	1-utt_wo_finetune	1-utt_w_finetune	50-utts_w_finetune	N/P	p
1-utt_wo_finetune vs 50-utts_w_finetune	11.75	–	66.00	22.25	<0.01
1-utt_w_finetune vs 50-utts_w_finetune	–	19.50	60.75	19.75	<0.01
1-utt_wo_finetune vs 1-utt_w_finetune	27.25	36.00	–	36.75	<0.01

V. CONCLUSION

This paper has proposed a context-aware mask prediction network for the end-to-end text-based speech editing task, which can delete, replace and insert the speech at the word level by editing the transcription. To simulate the speech editing process at the training stage, the text-based speech editing task is viewed as a two-stage process: masking and prediction, and a coarse-to-fine decoding method is proposed to achieve context-aware prediction. At the inference stage, three operations are designed based on CampNet, corresponding to the deletion, insertion, and replacement operations. Then, to synthesize the speech of arbitrary length text in insertion and replacement operations, a word-level autoregressive generation method is proposed. Finally, we propose a one-sentence speaker adaptation method for the CampNet and explore the ability of few-shot and one-shot learning based on CampNet, which can boost performance further by using only one sentence. Compared with TTS and VC, it also provides a new method for speech forgery. The experimental results demonstrate that the CampNet is better than the TTS, VoCo, and manual editing in subjective evaluation, objective evaluation, and operational efficiency for the text-based speech editing task. In addition, the few-shot learning ability based on CampNet is better than TTS and VC systems. Improving the speech quality further based on CampNet is the future work.

VI. ACKNOWLEDGMENT

This work is supported by the National Key Research and Development Plan of China (No.2020AAA0140003), the National Natural Science Foundation of China (NSFC) (No.62101553, No.61901473, No.61831022), the Key Research Project (No.2019KD0AD01), and is also partially funded by Huawei Noah’s Ark Lab.

REFERENCES

- [1] R. Derry, *PC audio editing with Adobe Audition 2.0: Broadcast, desktop and CD audio production*. CRC Press, 2012.
- [2] S. Whittaker and B. Amento, “Semantic speech editing,” in *Proceedings of the SIGCHI conference on Human factors in computing systems*, 2004, pp. 527–534.
- [3] Z. Jin *et al.*, “Speech synthesis for text-based editing of audio narration,” 2018.
- [4] M. Morrison, L. Rencker, Z. Jin, N. J. Bryan, J.-P. Caceres, and B. Pardo, “Context-aware prosody correction for text-based speech editing,” in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 7038–7042.
- [5] E. Moulines and F. Charpentier, “Pitch-synchronous waveform processing techniques for text-to-speech synthesis using diphones,” *Speech communication*, vol. 9, no. 5-6, pp. 453–467, 1990.
- [6] 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.
- [7] 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.
- [8] M. Airaksinen, L. Juvela, B. Bollepalli, J. Yamagishi, and P. Alku, “A comparison between straight, glottal, and sinusoidal vocoding in statistical parametric speech synthesis,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 26, no. 9, pp. 1658–1670, 2018.
- [9] Q. Hu, K. Richmond, J. Yamagishi, and J. Latorre, “An experimental comparison of multiple vocoder types,” in *Eighth ISCA Workshop on Speech Synthesis*, 2013.
- [10] A. v. d. Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, and K. Kavukcuoglu, “Wavenet: A generative model for raw audio,” *arXiv preprint arXiv:1609.03499*, 2016.
- [11] 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*. PMLR, 2018, pp. 2410–2419.
- [12] R. Prenger, R. Valle, and B. Catanzaro, “Waveglow: A flow-based generative network for speech synthesis,” in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 3617–3621.
- [13] R. Yamamoto, E. Song, and J.-M. Kim, “Parallel wavenet: A fast waveform generation model based on generative adversarial networks with multi-resolution spectrogram,” in *ICASSP 2020-*

- 2020 *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 6199–6203.
- [14] K. Kumar, R. Kumar, T. de Boissiere, L. Gestin, W. Z. Teoh, J. Sotelo, A. de Brébisson, Y. Bengio, and A. Courville, “Mel-gan: Generative adversarial networks for conditional waveform synthesis,” *arXiv preprint arXiv:1910.06711*, 2019.
- [15] J. Kong, J. Kim, and J. Bae, “Hifi-gan: Generative adversarial networks for efficient and high fidelity speech synthesis,” *arXiv preprint arXiv:2010.05646*, 2020.
- [16] 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)*. IEEE, 2019, pp. 5891–5895.
- [17] J. Su, Z. Jin, and A. Finkelstein, “Hifi-gan: High-fidelity denoising and dereverberation based on speech deep features in adversarial networks,” 2020.
- [18] 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.
- [19] J. Shen, R. Pang, R. J. Weiss, M. Schuster, N. Jaitly, Z. Yang, Z. Chen, Y. Zhang, Y. Wang, R. Skerry-Ryan *et al.*, “Natural tts synthesis by conditioning wavenet on mel spectrogram predictions,” in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 4779–4783.
- [20] Y. Wang, D. Stanton, Y. Zhang, R.-S. Ryan, E. Battenberg, J. Shor, Y. Xiao, Y. Jia, F. Ren, and R. A. Saurous, “Style tokens: Unsupervised style modeling, control and transfer in end-to-end speech synthesis,” in *International Conference on Machine Learning*. PMLR, 2018, pp. 5180–5189.
- [21] Y. Jia, Y. Zhang, R. J. Weiss, Q. Wang, J. Shen, F. Ren, Z. Chen, P. Nguyen, R. Pang, I. L. Moreno *et al.*, “Transfer learning from speaker verification to multispeaker text-to-speech synthesis,” *arXiv preprint arXiv:1806.04558*, 2018.
- [22] E. Nachmani, A. Polyak, Y. Taigman, and L. Wolf, “Fitting new speakers based on a short untranscribed sample,” in *International Conference on Machine Learning*. PMLR, 2018, pp. 3683–3691.
- [23] S.-X. Zhang, Z. Chen, Y. Zhao, J. Li, and Y. Gong, “End-to-end attention based text-dependent speaker verification,” in *2016 IEEE Spoken Language Technology Workshop (SLT)*. IEEE, 2016, pp. 171–178.
- [24] B. Sisman, J. Yamagishi, S. King, and H. Li, “An overview of voice conversion and its challenges: From statistical modeling to deep learning,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 132–157, 2021.
- [25] S. H. Mohammadi and A. Kain, “An overview of voice conversion systems,” *Speech Communication*, vol. 88, pp. 65–82, 2017.
- [26] Z. Jin, G. J. Mysore, S. Diverdi, J. Lu, and A. Finkelstein, “Voco: Text-based insertion and replacement in audio narration,” *ACM Transactions on Graphics (TOG)*, vol. 36, no. 4, pp. 1–13, 2017.
- [27] L. Sun, K. Li, H. Wang, S. Kang, and H. Meng, “Phonetic posteriorgrams for many-to-one voice conversion without parallel data training,” in *2016 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 2016, pp. 1–6.
- [28] 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,” *arXiv preprint arXiv:2107.01554*, 2021.
- [29] Z. Wu, O. Watts, and S. King, “Merlin: An open source neural network speech synthesis system,” in *SSW*, 2016, pp. 202–207.
- [30] N. Li, S. Liu, Y. Liu, S. Zhao, and M. Liu, “Neural speech synthesis with transformer network,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, 2019, pp. 6706–6713.
- [31] V. Nair and G. E. Hinton, “Rectified linear units improve restricted boltzmann machines,” in *Icml*, 2010.
- [32] C. Veaux, J. Yamagishi, K. MacDonald *et al.*, “Cstr vctk corpus: English multi-speaker corpus for cstr voice cloning toolkit,” *University of Edinburgh. The Centre for Speech Technology Research (CSTR)*, 2017.
- [33] 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.
- [34] T. Gulzar, A. Singh, and S. Sharma, “Comparative analysis of lpcc, mfcc and bfcc for the recognition of hindi words using artificial neural networks,” *International Journal of Computer Applications*, vol. 101, no. 12, pp. 22–27, 2014.
- [35] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [36] D. Povey, A. Ghoshal, G. Boulianne, L. Burget, O. Glembek, N. Goel, M. Hannemann, P. Motlicek, Y. Qian, P. Schwarz *et al.*, “The kaldi speech recognition toolkit,” in *IEEE 2011 workshop on automatic speech recognition and understanding*, no. CONF. IEEE Signal Processing Society, 2011.
- [37] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Advances in neural information processing systems*, 2017, pp. 5998–6008.
- [38] I. Goodfellow, Y. Bengio, and A. Courville, *Deep learning*. MIT press, 2016.
- [39] S. Ioffe and C. Szegedy, “Batch normalization: Accelerating deep network training by reducing internal covariate shift,” in *International conference on machine learning*. PMLR, 2015, pp. 448–456.
- [40] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, “Dropout: a simple way to prevent neural networks from overfitting,” *The journal of machine learning research*, vol. 15, no. 1, pp. 1929–1958, 2014.
- [41] R. Kubichek, “Mel-cepstral distance measure for objective speech quality assessment,” in *Proceedings of IEEE Pacific Rim Conference on Communications Computers and Signal Processing*, vol. 1. IEEE, 1993, pp. 125–128.
- [42] M. Müller, “Dynamic time warping,” *Information retrieval for music and motion*, pp. 69–84, 2007.
- [43] Y. Ai and Z.-H. Ling, “A neural vocoder with hierarchical generation of amplitude and phase spectra for statistical parametric speech synthesis,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 28, pp. 839–851, 2020.
- [44] A. Joshi, S. Kale, S. Chandel, and D. K. Pal, “Likert scale: Explored and explained,” *British journal of applied science & technology*, vol. 7, no. 4, p. 396, 2015.



Tao Wang received the B.E. degree from the Department of Control Science and Engineering, Shandong University (SDU), Jinan, China, in 2018. He is currently working toward the Ph.D. degree with the National Laboratory of Pattern Recognition, Institute of Automation (NLPR), Chinese Academy of Sciences (CASIA), Beijing, China. His current research interests include speech synthesis, voice conversion, machine learning, and transfer learning.



Jiangyan Yi received the Ph.D. degree from the University of Chinese Academy of Sciences, Beijing, China, in 2018, and the M.A. degree from the Graduate School of Chinese Academy of Social Sciences, Beijing, China, in 2010. She was a Senior R&D Engineer with Alibaba Group from 2011 to 2014. She is currently an Assistant Professor with the National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences, Beijing, China. Her current research interests include speech processing, speech recognition, distributed

computing, deep learning, and transfer learning.



Ruibo Fu is an assistant professor in National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences, Beijing. He obtained B.E. from Beijing University of Aeronautics and Astronautics in 2015 and Ph.D. from Institute of Automation, Chinese Academy of Sciences in 2020. His research interest is speech synthesis and transfer learning. He has published more than 10 papers in international conferences and journals such as ICASSP and INTERSPEECH and has won the best paper award twice in NCMMS 2017 and 2019.



Jianhua Tao (SM'10) received the Ph.D. degree from Tsinghua University, Beijing, China, in 2001, and the M.S. degree from Nanjing University, Nanjing, China, in 1996. He is currently a Professor with NLPR, Institute of Automation, Chinese Academy of Sciences, Beijing, China. His current research interests include speech synthesis and coding methods, human computer interaction, multimedia information processing, and pattern recognition. He has authored or coauthored more than 80 papers on major journals and proceedings including IEEE

TRANSACTIONS ON AUDIO, SPEECH, AND LANGUAGE PROCESSING, and received several awards from the important conferences, such as Eurospeech, NCMMS, etc. He serves as the chair or program committee member for several major conferences, including ICPR, ACII, ICMI, ISCSLP, NCMMS, etc. He also serves as the steering committee member for IEEE Transactions on Affective Computing, an Associate Editor for Journal on Multimodal User Interface and International Journal on Synthetic Emotions, the Deputy Editor-in-Chief for Chinese Journal of Phonetics.



Zhengqi Wen received the B.E. degree from the Department of Automation, University of Science and Technology of China, Hefei, China, in 2008 and the Ph.D. degree from the National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences, Beijing, China, in 2013. From March 2009 to June 2009, he was an intern student with Nokia Research Center, China. From December 2011 to March 2012, he was an intern student with the Faculty of Systems Engineering, Wakayama University, Japan. From July 2014 to

January 2015, he was a visiting scholar, under the supervision of Professor Chin-Hui Lee, with the School of Electrical and Computer Engineering, Georgia Institute of Technology, USA. He is currently an Associate Professor with the National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences, Beijing, China.