

Voice Cloning

Saiesh Prabhu Verlekar¹, Saili Kulkarni², Varad Naik², Aaron Mendes², Saiesh Naik²

¹Project Guide, Department of Information Technology, Shree Rayeshwar Institute of Engineering and Information Technology Goa, India

²Students, Department of Information Technology, Shree Rayeshwar Institute of Engineering and Information Technology Goa, India

Submitted: 25-01-2022

Revised: 05-02-2022

Accepted: 08-02-2022

ABSTRACT

Deep learning models for natural voice cloning methods begin in 2016, since then the main focus of the researchers is to make the voice more natural and get the output voice in real-time. Previously it used to take many hours of voice samples to clone a few seconds. After using deep learning models now, it is reduced to a few seconds. In this paper, we will study different techniques used for voice cloning. Some of these methods are multi-speaker generative model, speaker adoption, speaker encoding, vector quantization, etc. Index Terms - Voice cloning, text-to-speech, voice conversation, speaker adaptation, speaker encoding, vector quantization, neural network.

I. INTRODUCTION

The definition of VC (voice conversion) is to change the voice of the source speaker to the target speaker's voice without changing the language features. To imitate the target speaker, a VC system should modify the tone, accent, and pronunciation of the voice of the source speaker. The process of producing natural speech from a written input remains a challenging task despite decades of research [2]. These days there are many text-to-speech (TTS) methods able to get great results in terms of synthesis of natural voices very close to human once. Unfortunately, many of these systems learn to synthesize text only with a single voice. The goal of this project is to create a system that can effectively generate a data efficient natural speech for many different speakers. In practice, this technology can be used in many applications such as entertainment, creativity and this technology also allows restoring of voice or customizing digital assistants such as Siri [5].

Voice conversion consists of supervised and unsupervised scenarios. In the earlier times, researches on one-to-one or many to-one

conversion systems such as GMM-based and regression-based models, which convert one or many speakers voice into specific target speaker, is supervised learning. It learns to map the data into specific target distribution and usually achieves higher voice similarity and speech quality. But they need frame-level alignment on training data which may lack flexibility [10].

Unsupervised VC, which does not need parallel data, became a promising research direction to tackle the many-to-many voice conversion problem, and most of the works achieved good performance on seen speakers. Generative adversarial network (GAN) is one of the successful voice conversion methods owing to its distribution guarantee between the generated data and true data, and the difference between speakers is the most obvious one. Some works utilize vector quantization technique to extract content information and feed into WaveNet combined with one-hot speaker embedding to synthesize the voice. This output generated is remarkably close to natural human voice. However, they still suffer one problem which is that they cannot synthesize the voice which does not exist in the training data [2]. One-shot technique can solve the unseen speaker problem. One-shot model only needs one utterance from the source speaker and the target speaker. In the training phase, they should learn to differentiate the content information and speaker information, and most of them are not trained on speaker label to avoid overfitting on the training data. Proposed IN-based one-shot VC, which is a technique widely used in the computer vision, and it shows that this approach can learn the meaningful speaker embedding.

The important elements are:

Multi speaker generative model- In speech synthesis, generative models (speaker generative model) can be conditioned on text and speaker

identity. Although the text contains language proficiency and controls the content of the produced speech, the speaker identity captures features such as tone, speech rate, and pronunciation[1].

Speaker adaptation-
Speaker adaptation is based on fine-tuning a multi-speaker generative model with a few cloning samples[1].

Speaker encoding - Derives an embedding from the short utterance of a single speaker. The embedding is a meaningful representation of the voice of the speaker, such that similar voices are close in latent space [1].

Vector quantization - Vector quantization is a method of constraining an input from a continuous or otherwise large set of values by using a shared codebook [2].

II. LITERATURE SURVEY

According to Shijun Wang and Damian Borth they tried to perform voice conversion that could transform voice from audio source without losing the important contents in a language. They performed methods such as voice conversion, vector quantization, contrastive predictive coding, augmentation. The results that were found out were with contrastive predictive coding and Noise Augmentation the content encoder achieves very good results in capturing the content information. Meanwhile they also concluded that speaker encoder was successfully able to extract speaker characteristics[7].

Ju-chieh Chou, Cheng-chieh Yeh, Hung-yi Lee proposed a model that could effectively perform one shot voice conversion with a single utterance from the source and the speaker, before this voice conversion with parallel data was not successful. They used disentangled representations approach, instance normalization and generative model, using this method and model they were able to perform voice conversion on unseen speakers as well as the subjective and objective progression showed positive results for the target speaker[3].

Hieu-Thi Luong and Junichii Yamagishi

$$\min_{W, e} \mathbb{E}_{\substack{s_i \sim S_i \\ (t_{i,j}, a_{i,j}) \sim T_{s_i}}} \{L(f(t_{i,j}, s_i; W, e_{s_i}), a_{i,j})\}$$

When S means a set of speakers, T_{s_i} means a training set of text-audio pairs for speaker s_i , $a_{i,j}$ means the ground-truth audio for $t_{i,j}$ of speaker s_i . The expectation is guessed over text-audio pairs of each training speakers. \square and \square are used to designate the trained

Nautilus: A Versatile Voice Cloning System It can act as a TTS or VC system with high consistency in terms of speaker characteristics when switched between two[8].

Kaushik Daspute, Hemang Pandit Real Time Voice Cloning In this they tried to replicate three phase pipe line model that allows you to clone unseen voice with few sec voice sample[9].

Sercan Ö. Arık, Jitong Chen, Kainan Peng, Weiping, Yanqi Zhou. Neural Voice Cloning with a Few Samples In this they discussed about speaker encoding and speaker adoption, speaker encoding gives fast output but the voice feels little synthetic as compared to speaker adoption [1].

Merlijn Blaauw, Jordi Bonada and Ryunosuke Daido DATA EFFICIENT VOICE CLONING FOR NEURAL SINGING SYNTHESIS From small amount of voice sample we can produce voice fitting or speaker adoption [6].

Paarth Neekhara, Shehzeen Hussain, Shlomo Dubnov Expressive Neural Voice Cloning Evaluating 3 aspects of extracting and transferring audio reference to speech [4].

Da-Yi Wu, Hung-yi Lee proposed a vector quantization based one-shot VC approach without any supervision on speaker labels. Their experiments showed that the VQVC learns a meaningful embedding space without any supervision. Then they even performed VC to unseen speakers with only one utterance, and subjective evaluations showed good results in terms of similar to target speakers[2].

III. ALGORITHM

1.1 Multi-Speaker Generative Modeling to Voice Cloning

Consider multi-speaker generative model, $f(t_{i,j}, s_i; W, e_{s_i})$, which takes text input $t_{i,j}$ and a speaker samples s_i . Trainable parameters are parameterized by W , and e_{s_i} . The latter s_i corresponding to the trainable speaker embedding. W and e_{s_i} are optimized by minimizing L (loss function) that penalizes the difference between generated and ground-truth audios (eg.: a regression loss for spectrograms):

parameters and embeddings. The speaker embeddings have been shown to effectively capture speaker characteristics with low-dimensional vectors. Generative loss, discriminatory features (e.g., gender and verbal communication) are seen in the speaker embedding space.

For cloning voice, we take out the speaker characteristics of an unknown speaker s_k from a set of cloning audios A_{sk} , and generate an audio given any text for that speaker. Speech naturalness and speaker similarity are the two-performance metrics for the generated audio [1].

1.2 Speaker adaptation

The idea is to fine-tune a trained multi-speaker model for an unseen speaker using audio-text pairs. It can be applied to either the speaker embedding or the entire model. For embedding-only adaptation, we have the following:

$$\min_{e_{s_k}} \mathbb{E}_{(t_{k,j}, a_{k,j}) \sim T_{s_k}} \left\{ L \left(f(t_{k,j}, s_k; \widehat{W}, e_{s_k}) \right) \right\}$$

Where T_{sk} is a set of text-audio pairs for the target

$$\min_{W, \Theta} \mathbb{E}_{s_i \sim S, (t_{i,j}, a_{i,j}) \sim T_{s_i}} \left\{ L \left(f(t_{i,j}, s_i; W, g(\mathcal{A}_{s_i}; \Theta)), a_{i,j} \right) \right\}.$$

$$\min_{W, e_{s_k}} \mathbb{E}_{(t_{k,j}, a_{k,j}) \sim T_{s_k}} \left\{ L \left(f(t_{k,j}, s_k; W, e_{s_k}), a_{k,j} \right) \right\}.$$

1.4 Vector Quantization

Vector quantization is a method of constraining an input from a continuous or otherwise large set of values by using a shared codebook. It also helps the decoder to do reconstruction with the quantized vectors. In signal domain let $X = \{x_1, x_2, \dots, x_T\}$ where T denotes the sequence of acoustic features. An encoder is applied to transform X into a sequence of representation $E = \{e_1, e_2, \dots, e_T\}$. Each value of the E set is quantized into a sequence of codes $Q = \{q_1, q_2, \dots, q_T\}$.

The quantization function can be written as :

$$q_t = \arg \min_{q \in Q} (||e_t - q||_4)$$

the function takes the encoded e_t and selects the closest q from the codebook based on similar distances. A decoder is employed to reconstruct the correct input features which might have been lost with input Q . The decoder gives the output in the form of $\hat{X} = \{\hat{x}_1, \hat{x}_2, \dots, \hat{x}_t\}$ [2].

1.5 Contrastive predictive coding

Contrastive predictive coding learns self-supervised representations by predicting the future in the latent space. This model encourages the model to capture time-variant information while discarding time-invariant features. Contrastive

speaker s_k . For whole model adaptation, we have the following:

The entire model provides more degrees of freedom for speaker adaptation but, its optimization is difficult for a small amount of cloning data. To avoid overfitting, early stopping is required [1].

1.3 Speaker encoding

In this case they suggested a speaker encoding method to directly measure the embedding of the speaker from the audio samples of the unseen speaker. Such type of model doesn't require any fine-tuning during voice cloning. so, the same model can be used for all unseen speakers. The speaker encoder, $g(A_{sk}; \Theta)$, takes a set of cloning audio samples A_{sk} and measures e_{sk} for speaker s_k . The model is parametrized by Θ [1]:

predictive coding loss is contrastive. With the contrastive loss, the prediction is achieved by minimizing the dot product provided by an encoder. In CPC the sequence from the encoder is passed to current network to produce context.

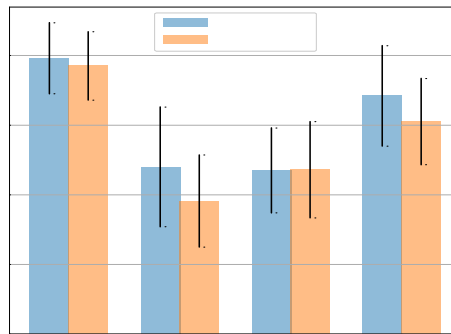
$C = \{ c_1, c_2, \dots, c_4 \}$. The main function of the model is to predict the future patterns given the current input [3], [5].

IV. OBSERVATION

Similarity test were conducted with Mean Opinion Score. In this test subjects were presented

with pairs of utterances. Subjects were asked to assign a score ranging from 1-5. 5) Same, absolutely sure 4) Slightly sure 3) Not sure 2) Slightly different 1) Different absolutely sure. They randomly selected 10 speakers from seen speakers and unseen speakers. MOS scores are presented in the below figure.

Blue (seen speaker), Pitch (unseen speaker)



Thus, we can show that both speaker adaptation and speaker encoding can achieve an MOS similar to the baseline. We also observe that drawbacks of training a multi-speaker model using a dataset of low-quality audios and limited speaker diversity. This test also showed us to complete VC to unseen speakers in one shot try [3],[5].

V. CONCLUSION

We have studied the different approaches for voice cloning like Speaker adaptation, speaker encoding, vector quantization, multi-speaker generative model, etc. From the learnings we found that both speaker adaptation and speaker encoding approaches can get good cloning quality even with a few audio sample, for naturalness both speaker adaptation & speaker encoding can achieve good results. On other hand, Speaker adaptation takes more time to give results but the quality of cloned audio is better than speaker encoding.

REFERENCE

- [1] Sercan Ö. Arik, Jitong Chen, Kainan Peng, Wei Ping, Yanqi Zhou. Neural Voice Cloning with a Few Samples. Baidu Research 1195 Bordeaux Dr. Sunnyvale, CA 94089.
- [2] Da-Yi Wu, Hung-yi Lee. One-shot Voice Conversion by Vector Quantization. College of Electrical Engineering and Computer Science, National Taiwan University (x07922119, hungvilee@ntu.edu.tw).
- [3] Ju-chieh Chou, Cheng-Chieh Yeh. To implement one shot voice cloning by separating speaker from source and target speaker.
- [4] Paarth neekhara, Shehzeen Hussain. Expressive neural voice cloning.
- [5] Giuseppe Ruggiero, Enrico Zovato. Voice cloning: A multi speaker TTS synthesis approach based on transfer learning.
- [6] Merlin Blaauw, Jordi Bonada, and Ryunosuke Daido. Data Efficient Voice Cloning For Neural Singing Synthesis.
- [7] Shilun Wang, Damian Borth. Towards High Quality Zero-Shot Voice Conversion.
- [8] Hieu-Thi Luong and Junichi Yamagishi. Nautilus: A Versatile Voice Cloning System.
- [9] Kaushik Daspute, Hemang Pandit. Real Time Voice Cloning.
- [10] Yen-Hao Chen, Da-vi wu. Again-VC A one shot voice conversion using activation guidance and AdalN.