One-shot Voice Conversion by Separating Speaker and Content Representations with Instance Normalization

Ju-Chieh Chou, Hung-yi Lee, Interspeech 2019.



Outline

- 1. Introduction
- 2. Proposed Approach
 - Model
 - Experiments
- 3. Conclusion

Outline

1. Introduction

- 2. Proposed Approach
 - Model
 - Experiments
- 3. Conclusion

Voice conversion

- Change the **characteristic** of an utterance while maintaining the language content the same.
- Characteristic: accent, speaker identity, emotion...
- This work: focuses on speaker identity conversion.



Conventional: supervised VC with parallel data

- Same sentences, different signal from 2 speakers.
- Formulated as a supervised learning problem.
- Problem: require parallel data, which is hard to collect.



Recently: unsupervised VC with non-parallel data

- Trained on non-parallel corpus, which is more attainable.
- Prior work: utilize deep generative model, ex. VAE, GAN, cycleGAN.
- Problem: cannot convert to speakers not in the training data.
- Our goal: train a model which is able to convert to speakers not in the training data.



Motivation

- Intuition: speech signals inherently carry both content and speaker information.
- Learn the content/speaker representation separately.
- Synthesize the target voice by combining the source content representation and target speaker representation.



Outline

- 1. Introduction
- 2. Proposed Approach
 - Model
 - Experiments
- 3. Conclusion

Model overview

- One-shot VC: use a utterance from target speaker as reference, and synthesize this reference speaker's voice.
- Idea: separately encode speaker and content information with some special designed layers.



Idea

IN

- Speaker information invariant within an utterance.
- **Content information** varying within an utterance.

Special Designed Layers:

Feature map Channel Instance Normalization Layer: normalizing speaker information (μ , σ) while preserving content information. M[']_c = σ_c

t=1

Intuition: normalize global information out (ex. high frequency), retain Charges Averation Pooling Layer: calculating speaker information (γ, β) . $M'_{c} = \sum_{r} \frac{M^{t}_{c}}{T}$

Adaptive Instance Normalization Layer: provide speaker, $= \gamma_c \frac{M_c - \mu_c}{\tau} + \beta_c$ Adall information (γ, β) .

Model - training

Problem: how to factorize the representations?



AVG calculating speaker information(γ , β).

IN normalizing speaker information (μ , σ) while preserving content information. AdaIN provide speaker information (γ , β).

Model - testing



AVG calculating speaker information(γ , β).

IN normalizing speaker information (μ , σ) while preserving content information. AdalN provide speaker information (γ , β).

Experiments – effect of IN

- Train another speaker classifier to see how much speaker information in content representations.
- The lower the accuracy is, the less speaker information it contains.
- Content encoder + IN: less speaker information.



Experiments – speaker embedding visualization

- Does speaker encoder learns meaningful representations?
- One color represents one speaker's utterances.

Unseen speakers'

utterances

• z_s from different speakers are well separated.

Speaker

Encoder E_s

AVG



Experiments - subjective

 Ask subjects to score the similarity between 2 utterances in 4-scales.





Experiments - subjective

Same, absolute sure Same, not sure



Different, absolute sure

- Ask subjects to score the similarity between 2 utterances in 4-scales.
- Our model is able to ulletgenerate the voice similar to target speaker's.



Demo page: https://jjery2243542.github.io/one-shotvc-demo/

Demo (unseen) Male to Male

Source: Target: Converted:

Female to Male

Source: Target: Target: Converted:

Conclusion

- We proposed a one-shot VC model, which is able to convert to unseen speaker with one reference utterance.
- By IN and AdaIN, our model is able to learn factorized representations.

Thank you for your attention.