

AutoVC: Zero-Shot Voice Style Transfer with Only Autoencoder Loss

Kaizhi Qian*¹, Yang Zhang*²³, Shiyu Chang²³,
Xuesong Yang¹, Mark Hasegawa-Johnson¹

¹University of Illinois at Urbana-Champaign

²MIT-IBM Watson AI Lab

³IBM Research Cambridge



Motivation

- Voice conversion aims to modify the source speech to make it sound like being uttered by another speaker.



IBM®

MIT

Motivation

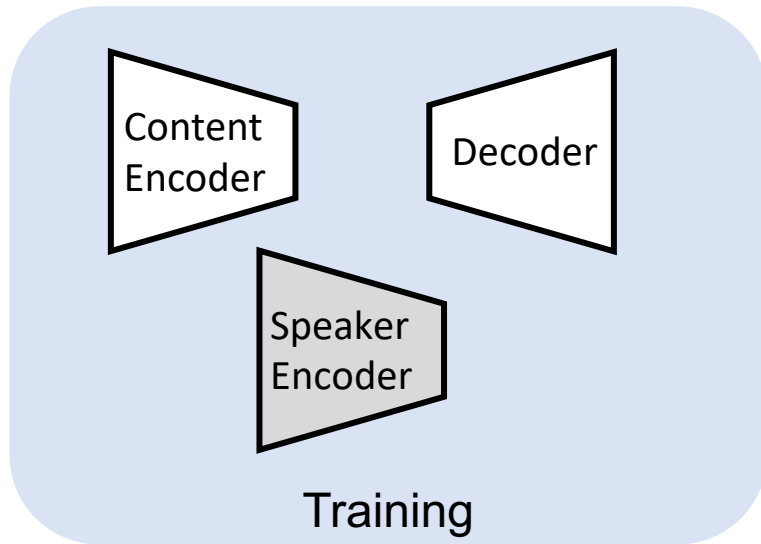
- Voice conversion aims to modify the source speech to make it sound like being uttered by another speaker.
- Existing voice style transfer techniques:
 - Use complex architectures and training schemes but do not work well for speech
 - Only convert between seen speakers



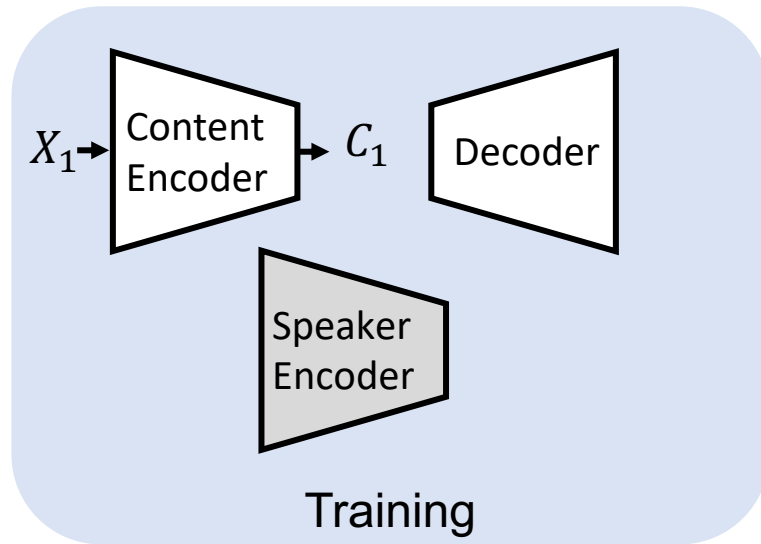
IBM®

MIT

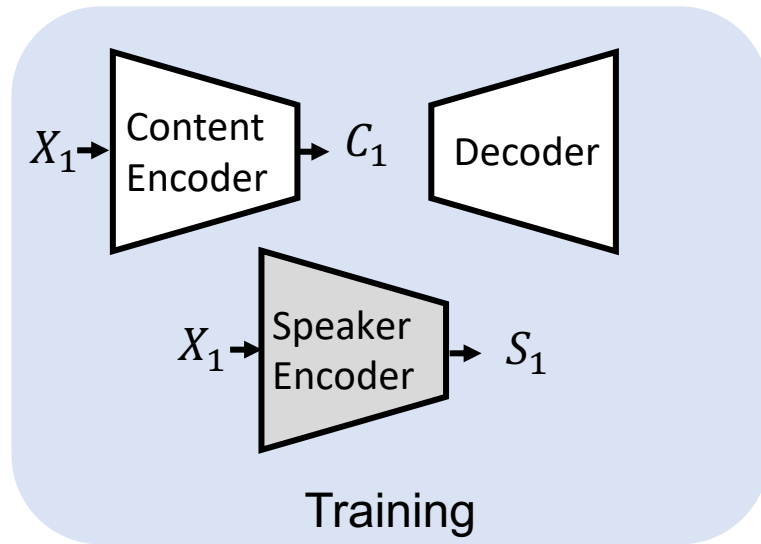
AutoVC



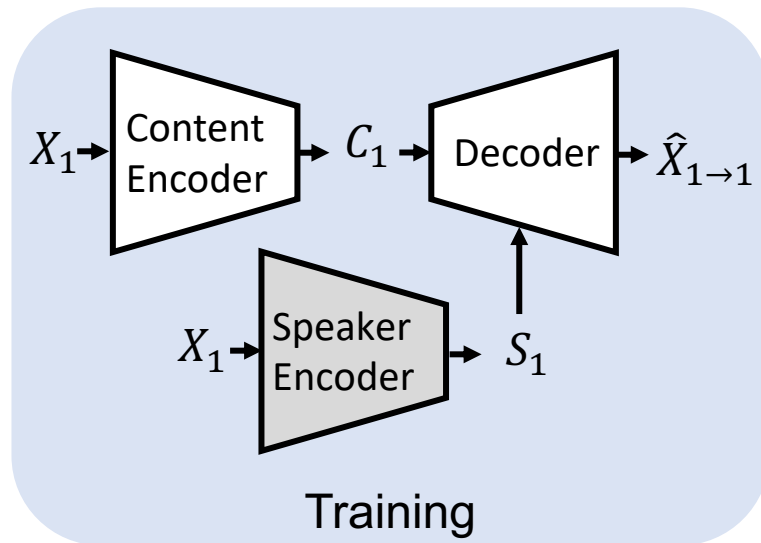
AutoVC



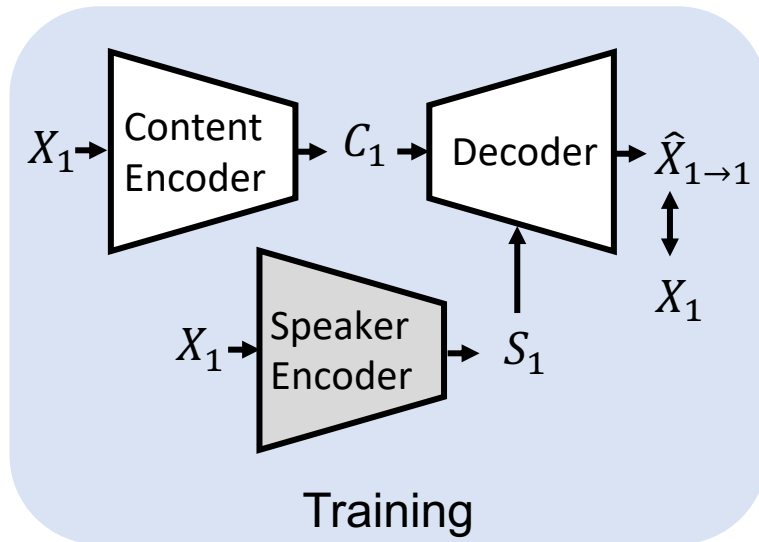
AutoVC



AutoVC



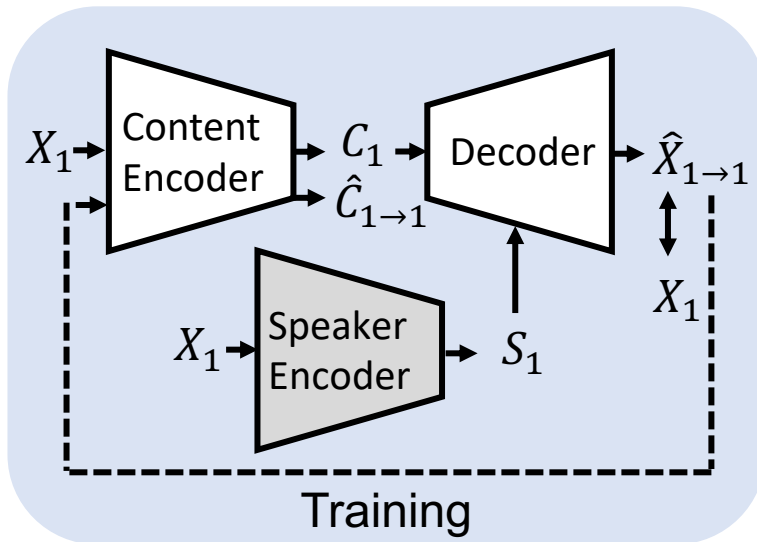
AutoVC



- Train only on **self-reconstruction** Loss:

$$\mathbb{E} \left[\|\hat{X}_{1 \rightarrow 1} - X_1\|_2^2 + \lambda \|\hat{C}_{1 \rightarrow 1} - C_1\|_1 \right]$$

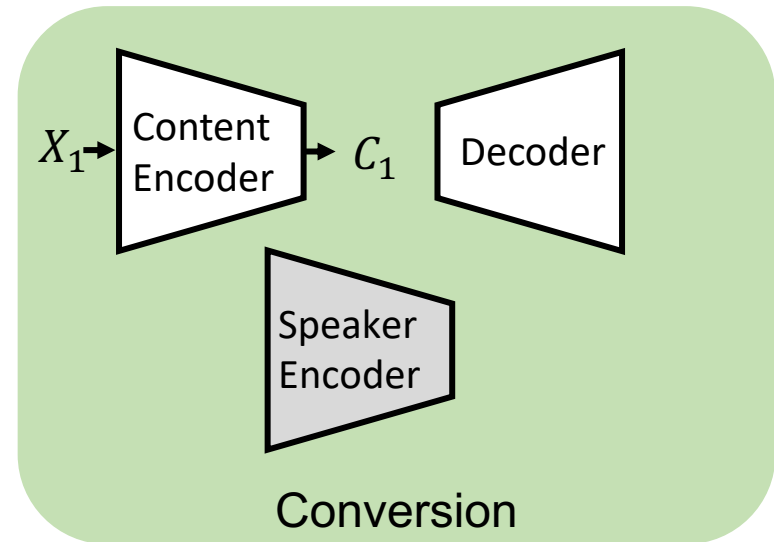
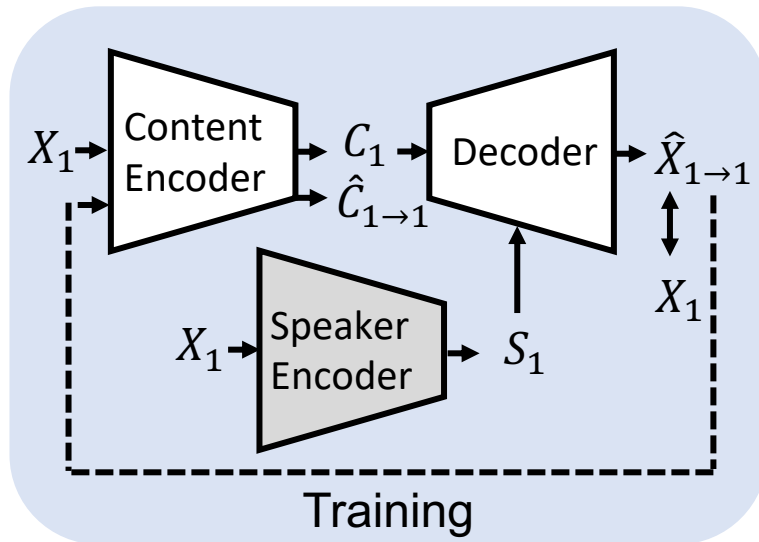
AutoVC



- Train only on **self-reconstruction** Loss:

$$\mathbb{E} \left[\|\hat{X}_{1 \rightarrow 1} - X_1\|_2^2 + \lambda \|\hat{C}_{1 \rightarrow 1} - C_1\|_1 \right]$$

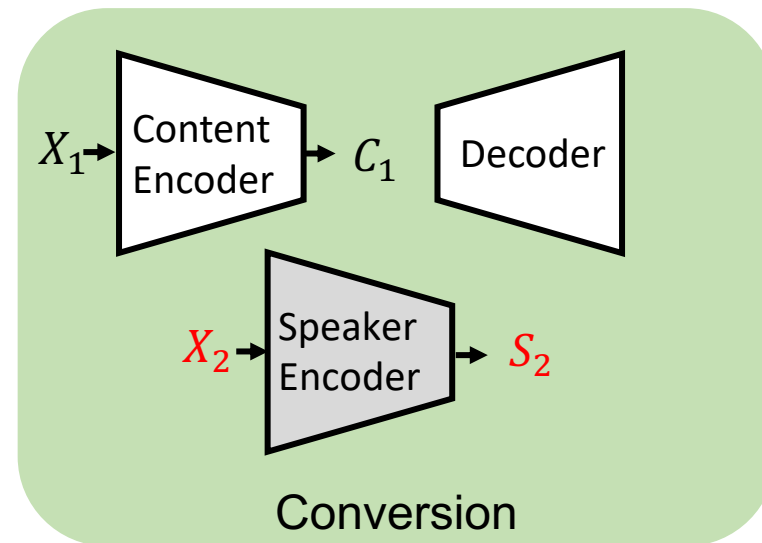
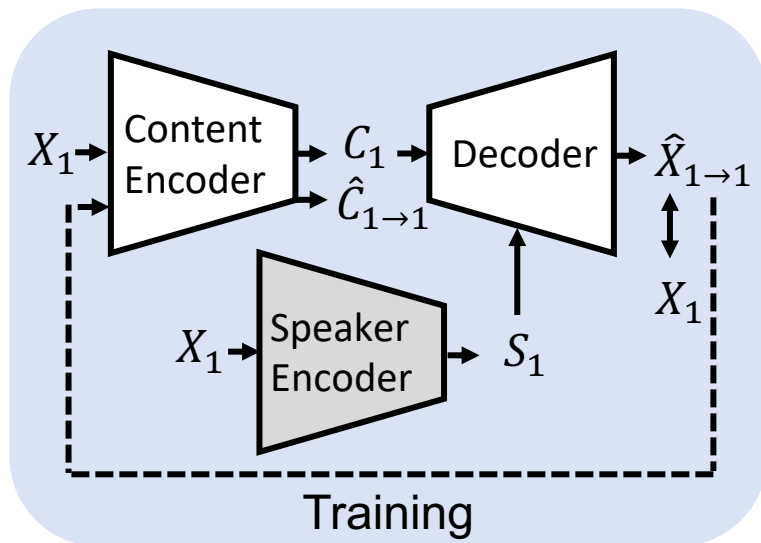
AutoVC



- Train only on **self-reconstruction** Loss:

$$\mathbb{E} \left[\|\hat{X}_{1 \rightarrow 1} - X_1\|_2^2 + \lambda \|\hat{C}_{1 \rightarrow 1} - C_1\|_1 \right]$$

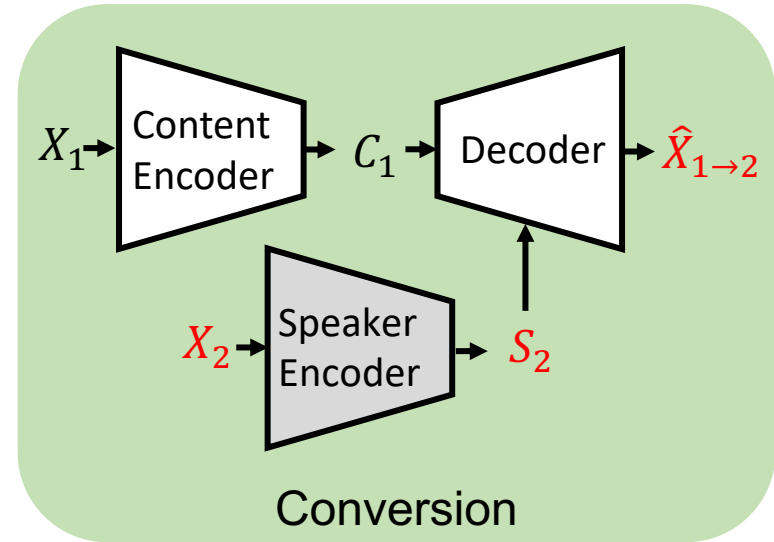
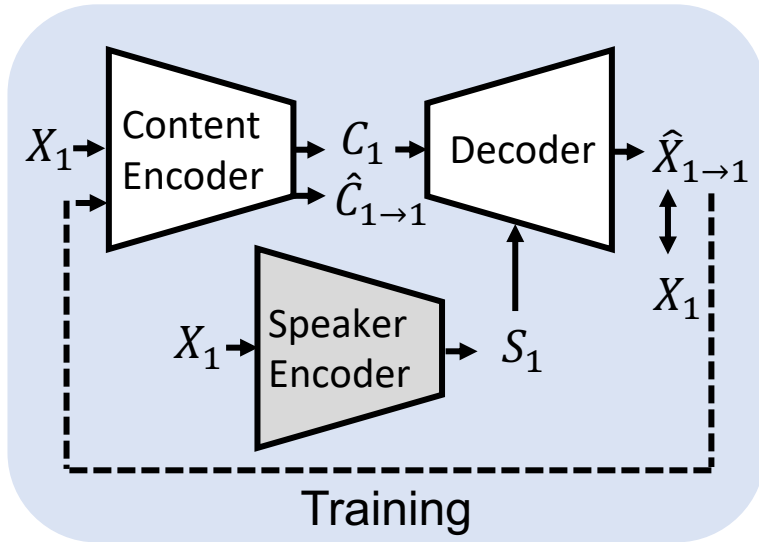
AutoVC



- Train only on **self-reconstruction** Loss:

$$\mathbb{E} \left[\|\hat{X}_{1 \rightarrow 1} - X_1\|_2^2 + \lambda \|\hat{C}_{1 \rightarrow 1} - C_1\|_1 \right]$$

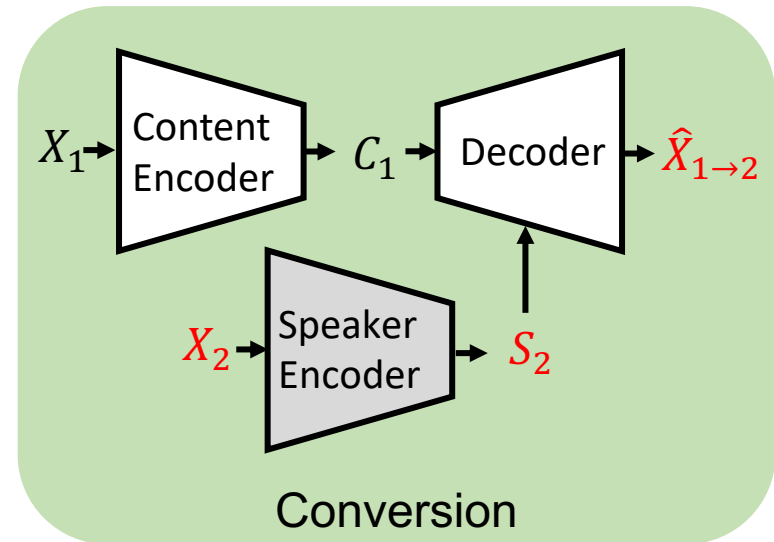
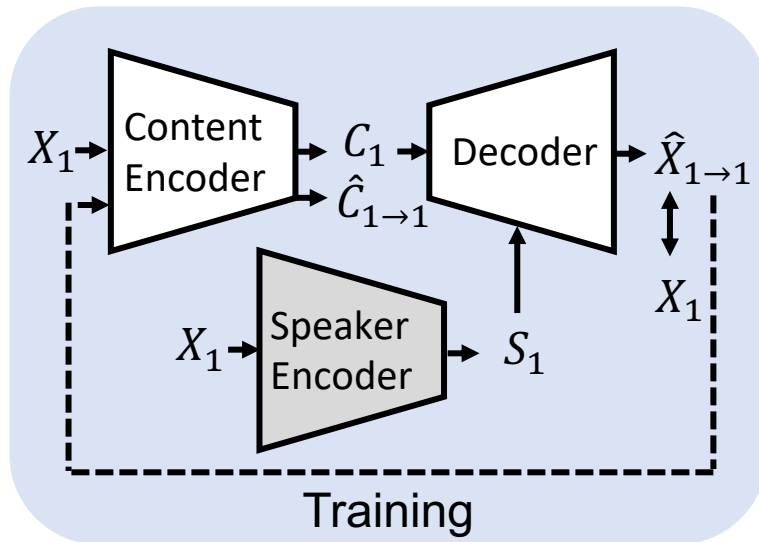
AutoVC



- Train only on **self-reconstruction** Loss:

$$\mathbb{E} \left[\|\hat{X}_{1 \rightarrow 1} - X_1\|_2^2 + \lambda \|\hat{C}_{1 \rightarrow 1} - C_1\|_1 \right]$$

AutoVC



- Train only on **self-reconstruction** Loss:

$$\mathbb{E} \left[\|\hat{X}_{1 \rightarrow 1} - X_1\|_2^2 + \lambda \|\hat{C}_{1 \rightarrow 1} - C_1\|_1 \right]$$

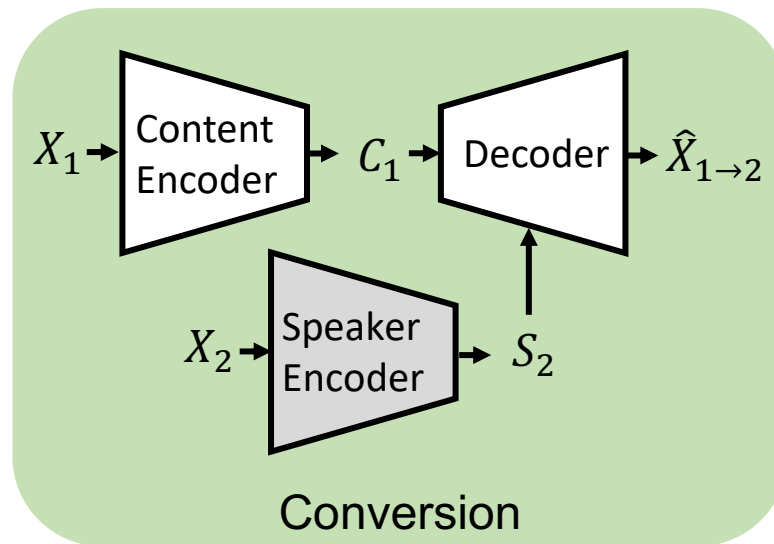
- With **bottleneck tuning**, AutoVC can **match the distribution!**



IBM

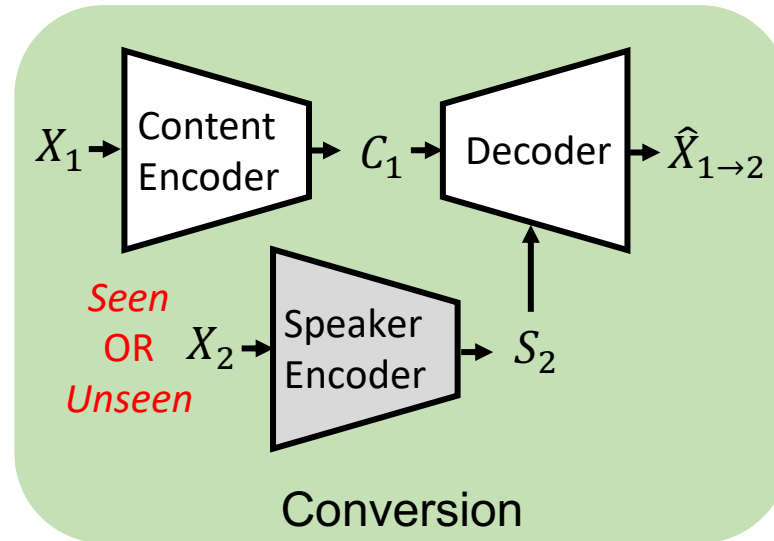


AutoVC



- Speaker encoder is pretrained

AutoVC



- Speaker encoder is pretrained
- Can generalize to **unseen speakers** – zero-shot conversion

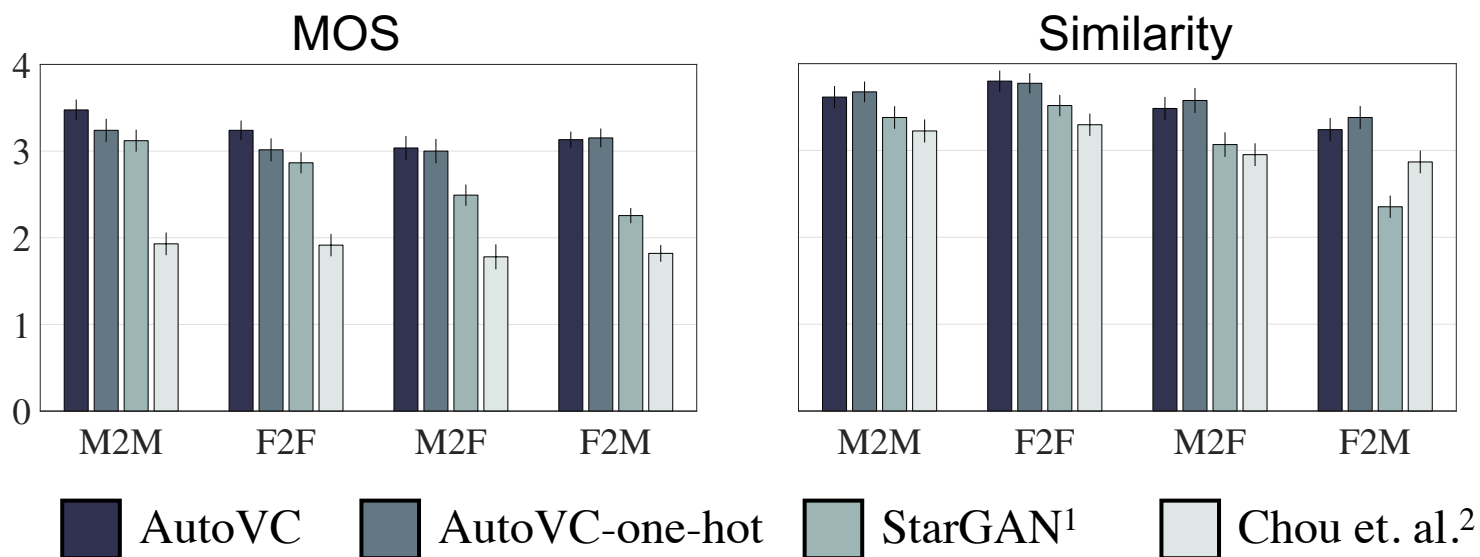


Conversion Between Seen Speakers

Source

Target

Converted



¹StarGAN-VC: Non-parallel many-to-many voice conversion with star generative adversarial networks

²Multi-target Voice Conversion without Parallel Data by Adversarially Learning Disentangled Audio Representations



Conversion Between Unseen Speakers

- The first zero-shot voice conversion framework

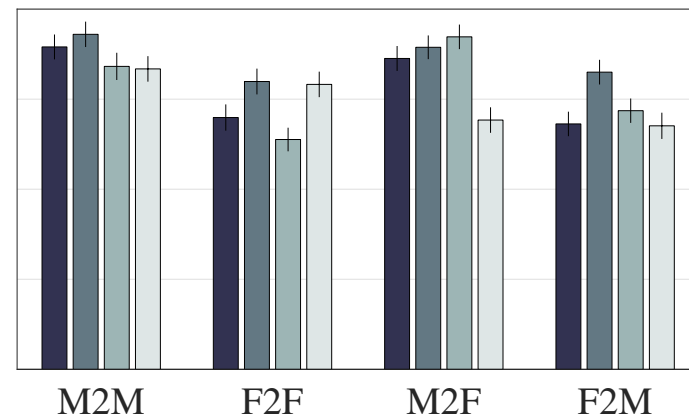
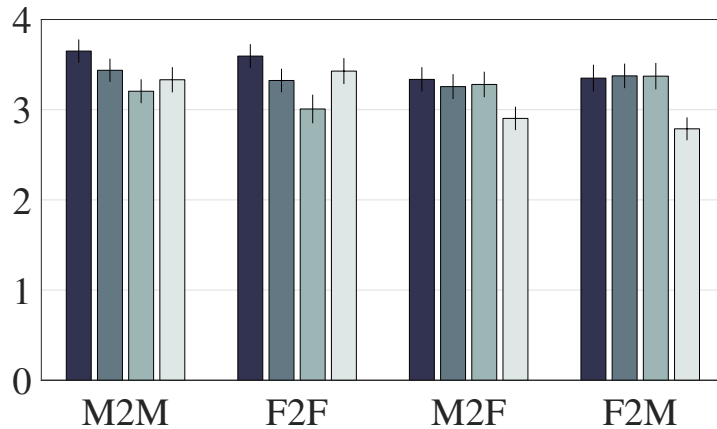
Source

Target

Converted

MOS

Similarity



■ Seen to seen ■ Seen to unseen ■ Unseen to seen ■ Unseen to unseen



Take Away

- Autoencoder is all you need to achieve theoretically ideal voice conversion



IBM®

MIT

Take Away

- Autoencoder is all you need to achieve theoretically ideal voice conversion
- AutoVC generalizes well to unseen speakers



IBM®

MIT

Thank you!

Poster #225