# *ACVAE-VC: Non-Parallel Voice Conversion With Auxiliary Classifier Variational Autoencoder* by Kameoka et al.

Master MVA - Audio Signal Analysis, Indexing, and Transformations

Dorian Desblancs

March 15, 2021

école
normale
supérieure
paris—saclay

université
PARIS-SACLAY

# Outline

1. Motivation

2. Methodology

3. Results

4. Wrap-up

# Motivation

# Voice Conversion

- Change one or more aspects of a speech signal while preserving linguistic information
- Applications: speaker-identity modification, speaking assistance, speech enhancement...
- Two types:
  - Parallel voice conversion
  - **Non-Parallel voice conversion**

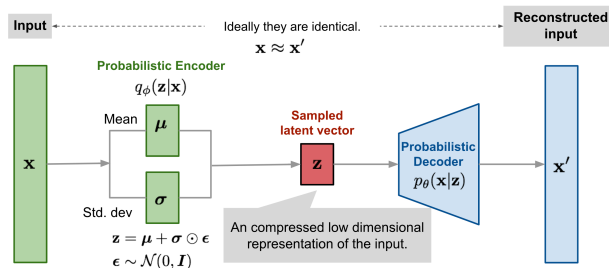Illustration of the variational autoencoder model with the multivariate Gaussian assumption (2) (3)

# Methodology

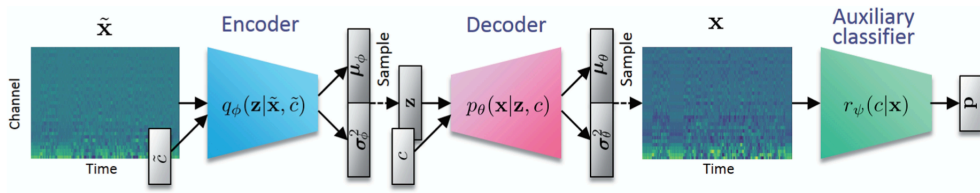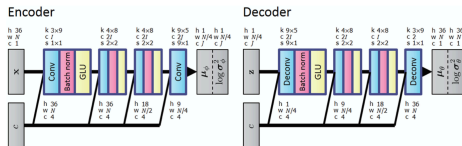Illustration of the auxiliary class variational autoencoder model (1)

VAE Architecture (1)



Auxiliary Model Architecture (1)

- Mel-Cepstral Coefficients



Sample Inputs and Outputs (1)

- Loss function:

$$L(\phi, \theta, \psi) = VAE(\phi, \theta) + \lambda_Q Q(\phi, \theta, \psi) + \lambda_R R(\psi) \tag{1}$$

- Performance Measure: mel-cepstral distortion (MCD) along DTW Path

$$MCD[dB](x, y) = \frac{10}{ln10} \sqrt{2 \sum_{d=2} D(x_d - y_d)^2} \tag{2}$$

- Data set: Voice Conversion Challenge (VCC) 2018

# Results

| Speakers | | Auxiliary classifier | |
| --- | --- | --- | --- |
| source | target | not included | included |
| SF1 | SM1 | $7.48 \pm 0.150$ | $\mathbf{6.70 \pm 0.129}$ |
| | SF2 | $7.38 \pm 0.163$ | $\mathbf{6.57 \pm 0.134}$ |
| | SM2 | $7.70 \pm 0.140$ | $\mathbf{6.97 \pm 0.124}$ |
| SM1 | SF1 | $7.64 \pm 0.144$ | $\mathbf{7.01 \pm 0.108}$ |
| | SF2 | $6.93 \pm 0.148$ | $\mathbf{6.29 \pm 0.133}$ |
| | SM2 | $7.25 \pm 0.136$ | $\mathbf{6.64 \pm 0.111}$ |
| SF2 | SF1 | $7.83 \pm 0.164$ | $\mathbf{6.94 \pm 0.115}$ |
| | SM1 | $7.25 \pm 0.151$ | $\mathbf{6.36 \pm 0.108}$ |
| | SM2 | $7.49 \pm 0.167$ | $\mathbf{6.85 \pm 0.137}$ |
| SM2 | SF1 | $7.82 \pm 0.176$ | $\mathbf{7.24 \pm 0.151}$ |
| | SM1 | $7.22 \pm 0.150$ | $\mathbf{6.66 \pm 0.133}$ |
| | SF2 | $7.15 \pm 0.170$ | $\mathbf{6.64 \pm 0.152}$ |

Results obtained with and without auxiliary classifier (1)

| Speakers | | non-parallel methods | | | | parallel method |
|---|---|---|---|---|---|---|
| source | target | VAE [19] | VAEGAN [20] | StarGAN [35] | Proposed | sprocket [61] |
| SF1 | SM1 | $7.66 \pm 0.123$ | $7.70 \pm 0.122$ | $7.81 \pm 0.126$ | $\mathbf{6.70 \pm 0.129}$ | $6.91 \pm 0.119$ |
| | SF2 | $7.53 \pm 0.118$ | $7.43 \pm 0.124$ | $7.54 \pm 0.146$ | $\mathbf{6.57 \pm 0.134}$ | $6.70 \pm 0.125$ |
| | SM2 | $8.06 \pm 0.143$ | $8.04 \pm 0.145$ | $8.11 \pm 0.123$ | $\mathbf{6.97 \pm 0.124}$ | $7.06 \pm 0.118$ |
| SM1 | SF1 | $8.25 \pm 0.104$ | $8.20 \pm 0.128$ | $8.27 \pm 0.119$ | $\mathbf{7.01 \pm 0.108}$ | $\mathbf{7.01 \pm 0.114}$ |
| | SF2 | $7.43 \pm 0.111$ | $7.23 \pm 0.117$ | $7.27 \pm 0.134$ | $\mathbf{6.29 \pm 0.133}$ | $6.30 \pm 0.108$ |
| | SM2 | $7.92 \pm 0.106$ | $7.82 \pm 0.103$ | $7.56 \pm 0.106$ | $6.64 \pm 0.111$ | $\mathbf{6.58 \pm 0.099}$ |
| SF2 | SF1 | $7.97 \pm 0.127$ | $7.83 \pm 0.121$ | $7.99 \pm 0.144$ | $\mathbf{6.94 \pm 0.115}$ | $7.21 \pm 0.111$ |
| | SM1 | $7.38 \pm 0.108$ | $7.37 \pm 0.097$ | $7.28 \pm 0.112$ | $\mathbf{6.36 \pm 0.108}$ | $6.77 \pm 0.108$ |
| | SM2 | $7.92 \pm 0.122$ | $7.78 \pm 0.109$ | $7.75 \pm 0.124$ | $\mathbf{6.85 \pm 0.137}$ | $\mathbf{6.85 \pm 0.115}$ |
| SM2 | SF1 | $8.33 \pm 0.148$ | $8.20 \pm 0.158$ | $8.30 \pm 0.189$ | $\mathbf{7.24 \pm 0.151}$ | $7.31 \pm 0.116$ |
| | SM1 | $7.73 \pm 0.138$ | $7.66 \pm 0.142$ | $7.44 \pm 0.122$ | $\mathbf{6.66 \pm 0.133}$ | $6.88 \pm 0.114$ |
| | SF2 | $7.74 \pm 0.135$ | $7.65 \pm 0.137$ | $7.53 \pm 0.154$ | $\mathbf{6.64 \pm 0.152}$ | $6.78 \pm 0.146$ |

Results compared to other methods (1)

All code and results on my GitHub!

Username: d-dawg78, repo: MVA_ASAIT

# Wrap-up

Advantages:

- Great results
- Clear distinction between outputs depending on class

Improvement:

- Improve data processing and loading for large-scale training feasibility

- Auxiliary classifier added to a GAN ?
- Improve pre-processing
- Data augmentation to create more training samples out of small data sets

école
normale
supérieure
paris—saclay

université
PARIS-SACLAY

Dorian Desblancs · Master MVA - Audio Signal Analysis, Indexing, and Transformations | Wrap-up · 12/12

# THANK YOU

## QUESTIONS?

# Bibliography

1. Kameoka, H., Kaneko, T., Tanaka, K., & Hojo, N. (2019). ACVAE-VC: Non-parallel voice conversion with auxiliary classifier variational autoencoder. IEEE/ACM Transactions on Audio, Speech, and Language Processing, 27(9), 1432-1443.

2. Weng, L. (2018, August 12). From Autoencoder to Beta-VAE. GitHub. https://lilianweng.github.io/lil-log/2018/08/12/from-autoencoder-to-beta-vae.html#vae-variational-autoencoder

3. Kingma, D. P., & Welling, M. (2013). Auto-encoding variational bayes. arXiv preprint arXiv:1312.6114.