

# ACVAE-VC : Non-Parallel Voice Conversion With Auxiliary Classifier Variational Autoencoder (2018)

—  
*Audio Signal Analysis course final project*

Kamil Akesbi

Master MVA, France

15/03/2022

école —  
normale —  
supérieure —  
paris — saclay —

université  
PARIS-SACLAY

## ① Introduction

- Paper Motivation

- Existing approaches and drawbacks

## ② The main architecture : the Variational Auto-Encoder

## ③ Proposed Method

- Overall architecture

- About the auxiliary classifier

- Architecture details

## ④ Signal conversion

## ⑤ Experiments and results

- Experimental settings

- Paper results : comparison to conventional methods

- Implementation and VC speech generation

## ⑥ Conclusion

## Paper Topic

**Voice conversion** : converting some aspects of a speech signal without changing its linguistic information.

**Applications** : speaker identity modification - speaking assistance - speech enhancement

**Problem** : Many VC methods requires parallel source and target speeches which can be costly to collect and align.

⇒ **Goal** : develop a non-parallel VC method

- Existing approach to Non-Parallel VC : **Conditional VAE** [Hsu et al. 2016]

⇒ Extended version of VAE where the encoder and decoder take an attribute class  $c$  as input.

## Drawbacks :

- Don't capture time dependencies
- Over-smoothed outputs
- The encoder and decoder are free to ignore the attribute class label  $c$

## Solutions proposed in the paper :

- Fully convolutional architectures to capture time dependencies.
- Auxiliary classifier to ensure that the attribute class information is not lost.

## ① Introduction

- Paper Motivation

- Existing approaches and drawbacks

## ② The main architecture : the Variational Auto-Encoder

## ③ Proposed Method

- Overall architecture

- About the auxiliary classifier

- Architecture details

## ④ Signal conversion

## ⑤ Experiments and results

- Experimental settings

- Paper results : comparison to conventional methods

- Implementation and VC speech generation

## ⑥ Conclusion

- **Goal** :  $q_{\phi}(z|x) \iff p_{\theta}(z|x)$

- Conditional VAE modification :

$$q_{\Phi}(z|x, c)$$
$$p_{\theta}(x|z, c)$$

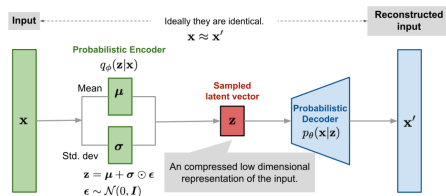


Figure 1: Variational Auto-Encoder

- **Lower bound** maximisation :

$$J(\phi, \theta) = \mathbb{E}_{(x,c) \sim p_d(x,c)} \left[ \mathbb{E}_{z \sim q(z|x,c)} [\log p(x|z, c)] \right] - \text{KL} [q(z|x, c) \parallel p(z)]$$

## ① Introduction

- Paper Motivation

- Existing approaches and drawbacks

## ② The main architecture : the Variational Auto-Encoder

## ③ Proposed Method

- Overall architecture

- About the auxiliary classifier

- Architecture details

## ④ Signal conversion

## ⑤ Experiments and results

- Experimental settings

- Paper results : comparison to conventional methods

- Implementation and VC speech generation

## ⑥ Conclusion

- VAE's Inputs and Outputs :

⇒ Sequence of 36 **Mel-Cepstral coefficients** [Tokuda et al. 1992]

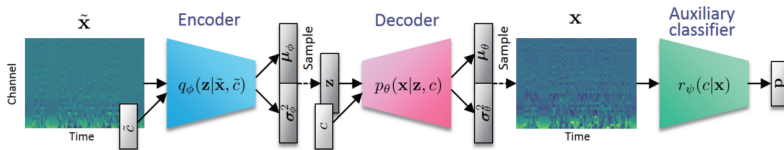


Figure 2: Structure of the proposed ACVAE-VC



- The auxiliary classifier helps to maximize the **mutual information** between the output  $x \sim p_\theta(x|z, c)$  and class attribute  $c|z$ .
- It can be seen as a **regularizer**.
- The final training criterion is given by :

$$\mathcal{J}(\phi, \theta) + \lambda_{\mathcal{Q}} \mathcal{Q}(\phi, \theta, \psi) + \lambda_{\mathcal{R}} \mathcal{R}(\psi)$$

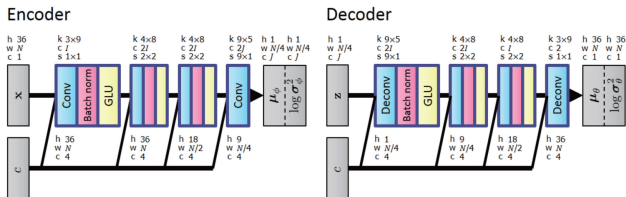


Figure 3: Encoder and Decoder architectures

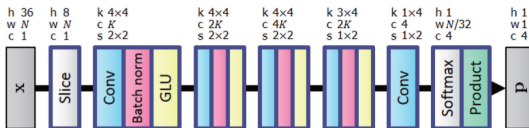


Figure 4: Auxiliary classifier architecture

## ① Introduction

- Paper Motivation

- Existing approaches and drawbacks

## ② The main architecture : the Variational Auto-Encoder

## ③ Proposed Method

- Overall architecture

- About the auxiliary classifier

- Architecture details

## ④ Signal conversion

## ⑤ Experiments and results

- Experimental settings

- Paper results : comparison to conventional methods

- Implementation and VC speech generation

## ⑥ Conclusion

**First** : Which output consider at test time ?

- Different possibilities :

$$\textcircled{1} \hat{x}_{mean} = \mu_{\theta}(\mu_{\phi}(x, c), \hat{c})$$

$$\textcircled{2} \hat{x}_{diff} = x - \bar{x}_{mean} + \hat{x}_{mean}$$

where  $\bar{x}_{mean} = \mu_{\theta}(\mu_{\phi}(x, c), c)$

$$\textcircled{3} \hat{x}_{samp} \sim p_{\theta}(x|\hat{z}, \hat{c})$$

**Second** : Reconstructing the speech signal using the **WORLD** vocoder  
[Morise et al. 2016]

## ① Introduction

- Paper Motivation

- Existing approaches and drawbacks

## ② The main architecture : the Variational Auto-Encoder

## ③ Proposed Method

- Overall architecture

- About the auxiliary classifier

- Architecture details

## ④ Signal conversion

## ⑤ Experiments and results

- Experimental settings

- Paper results : comparison to conventional methods

- Implementation and VC speech generation

## ⑥ Conclusion

- **Dataset** : Subset of the Voice Conversion Challenge 2018 dataset  
 $\Rightarrow$  Two female speakers, two male speakers
- **Performance measure** : average mel-cepstral distortion (MCDs) along the DTW path :

$$MCD[dB] = \frac{10}{\ln 10} \sqrt{2 \sum_{d=2}^D (x_d - y_d)^2}$$

Speakers		Layer type	
source	target	frame-independent	fully convolutional
SF1	SM1	$9.07 \pm 0.197$	<b><math>6.79 \pm 0.088</math></b>
	SF2	$8.73 \pm 0.121$	<b><math>6.51 \pm 0.096</math></b>
	SM2	$9.25 \pm 0.189$	<b><math>7.05 \pm 0.081</math></b>
SM1	SF1	$8.94 \pm 0.166$	<b><math>7.03 \pm 0.090</math></b>
	SF2	$8.33 \pm 0.214$	<b><math>6.29 \pm 0.095</math></b>
	SM2	$8.68 \pm 0.165$	<b><math>6.67 \pm 0.071</math></b>
SF2	SF1	$8.78 \pm 0.211$	<b><math>6.95 \pm 0.104</math></b>
	SM1	$8.54 \pm 0.198$	<b><math>6.45 \pm 0.103</math></b>
	SM2	$8.75 \pm 0.183$	<b><math>6.87 \pm 0.102</math></b>
SM2	SF1	$9.03 \pm 0.202$	<b><math>7.17 \pm 0.098</math></b>
	SM1	$8.65 \pm 0.182$	<b><math>6.63 \pm 0.083</math></b>
	SF2	$8.43 \pm 0.213$	<b><math>6.58 \pm 0.088</math></b>

Figure 5: Results obtained with or without taking into account time dependency

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

Figure 6: Results obtained with or without the auxiliary classifier

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

Figure 7: Comparison to the conventional non-parallel and parallel methods

- Code adapted from the public repository :  
<https://github.com/ariacat3366/ACVAE-VC>
- **Main modification** : save data preprocessing step to speed up training :  
⇒ More then 20 hours training to 30 mins
- Some hyper-parameters (sampling frequency,  $\lambda_Q$ ,  $\lambda_R$ ) + Reconstruction loss computation **differs from the paper...**  
→ Still qualitatively gives good results.
- Code and Results on my github :  
<https://github.com/kamilakesbi/MVA-Voice-Conversion-with-ACVAE>



## ① Introduction

- Paper Motivation

- Existing approaches and drawbacks

## ② The main architecture : the Variational Auto-Encoder

## ③ Proposed Method

- Overall architecture

- About the auxiliary classifier

- Architecture details

## ④ Signal conversion

## ⑤ Experiments and results

- Experimental settings

- Paper results : comparison to conventional methods

- Implementation and VC speech generation

## ⑥ Conclusion

## **Main contributions of the paper :**

This paper proposed a non-parallel VC method using a VAE variant called auxiliary classifier VAE (ACVAE)

- Key ideas :
  - ⇒ capturing time dependencies.
  - ⇒ Using information-theoretic regularization with the auxiliary classifier.

## **Results :**

- Better performances than conventional methods

## **Future work**

- Incorporating a neural vocoder instead of the WORLD vocoder.
- Testing Transformer approaches rather than the convolutional ones

- Hsu, Chin-Cheng et al. (2016). *Voice Conversion from Non-parallel Corpora Using Variational Auto-encoder*. URL: <https://arxiv.org/abs/1610.04019>.
- Morise, M. et al. (2016). *WORLD : A vocoder-based highquality speech synthesis system for real-time applications*. URL: [https://www.jstage.jst.go.jp/article/transinf/E99.D/7/E99.D\\_2015EDP7457/\\_article](https://www.jstage.jst.go.jp/article/transinf/E99.D/7/E99.D_2015EDP7457/_article).
- Tokuda, K. et al. (1992). *An adaptive algorithm for mel-cepstral analysis of speech*. URL: [https://www.researchgate.net/publication/3532054\\_An\\_adaptive\\_algorithm\\_for\\_mel-cepstral\\_analysis\\_of\\_speech](https://www.researchgate.net/publication/3532054_An_adaptive_algorithm_for_mel-cepstral_analysis_of_speech).

Thank you!