# Thesis

## Summaries on articles

## Article analisys

1. Voice conversion based on DNN for time-variant linear transformation

   - Linear transform is enough because we work with homo-domain mapping

<u>def</u>   1) Cepstrum — result of inv. fourier transform of the log of the estimated <u>signal spectrum</u>

2) Signal spectrum of a time series $x(t)$ — distribution of power into freq. components f, comp. that signal

- <u>GMM</u> based voice conversion approach:

(Gaussian mixture model)

joint proba of input and output features (on sepstrum space) is modeled by GMM; therefore, conversion input ↦ target is locally linear transformation

{ can be interp. as the use of the homo-domain condition }

- Another way — voice conv. based on non-negative matrix factorization

[not much was said about it]

- also <u>difference between speakers in cepstral space</u> was discussed:

a wave transformation via warp function
can be descr. as a linear transformation
in cepstral space.

$$\hat{z}^{-1} = \frac{z^{-1} - \alpha}{1 - \alpha z^{-1}}, \quad \text{where} \quad z = e^{j\omega} \qquad \alpha - \text{warping param.}$$
$$\hat{z} = e^{j\hat{\omega}} \qquad \omega - \text{initial,}$$
$$\hat{\omega} - \text{after transf.}$$

$\uparrow$ can be rewritten as

$$\hat{c} = AC, \quad \text{where} \quad A = \begin{pmatrix} 1-\alpha^2 & 2\alpha - 2\alpha^3 & \cdots \\ -\alpha + \alpha^3 & & \\ \cdots & & \end{pmatrix}$$

$\underbrace{\qquad\qquad\qquad\qquad}_{\text{rotation matr.}}$

- Deep NN for time variant linear transform.

$$x_t \longmapsto \{ \underset{\text{matrix}}{\text{lin. transform.}}, \text{Biases, } \underset{\text{param.}}{\text{warping}}, \underset{\text{feature}}{\text{dynamic}} \}$$

Linear transform: $\hat{y}_t = (A_t + A_t^{(\alpha)})(x_t - b_t^{src}) + b_t^{tgt}$

where $\qquad\qquad\qquad \Delta\hat{y}_t = G_{\Delta y}(x_t)$

$$G_{\Delta y}, G_A, G_\alpha, G_b \qquad A_t = G_A(x_t) \quad A_t^{(\alpha)} = \begin{pmatrix} 1-\alpha_t^2 \cdots \\ \cdots \end{pmatrix}$$

$\underbrace{\qquad\qquad\qquad}$

Sub networks $\qquad \alpha_t = G_\alpha(x_t)$
for parameters
estimation $\qquad [b_t^{src}, b_t^{tgt}]^T = G_b(x_t)$

## 2. Learning Latent Representations for Style control and transfer in end to end Speech synthesis.

- Main idea — learn latent repr. of the source audio to control style transfer

- Variational Autoencoder —
    reveals relationship betw. latent $z$ and observed $x$. True posterior $p_\theta(z|x)$ is impossible to find, therefore, likelihood $p_\theta(x)$ is indifferent. So $p_\theta(z|x)$ can be appoximated by $q_\phi(z|x)$ and

variational lower bound $\mathcal{L}(\theta, \phi; x)$:

$$\log p_\theta(x) = KL[q_\phi(z|x) \| p_\theta(z|x)] + \mathcal{L}(\theta, \phi; x)$$

$$\geq \mathcal{L}(\theta, \phi; x)$$

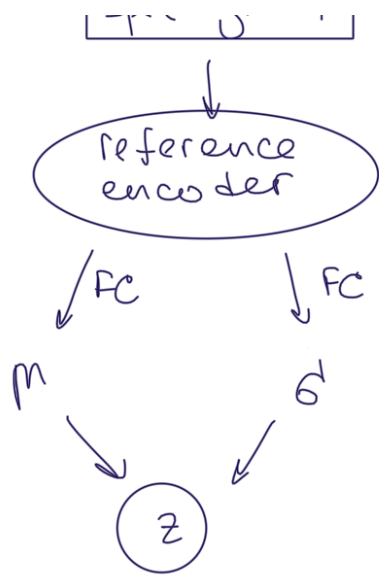$$= \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x|z)] - KL[q_\phi(z|x) \| p_\theta(z)]$$

where $p_\theta(z) \sim \mathcal{N}(z; 0, I)$, $q_\phi(z|x) \sim \mathcal{N}(z; \mu(x), \sigma^2(x) \cdot I)$

sort of decoder to decode latent $z$ to reconstruct $x$

Sampling $z$ from $\mathcal{N}(\mu, \sigma^2 \cdot I)$ is decomposed to sampling $\varepsilon \sim \mathcal{N}(0, I)$ and $z = \mu + \sigma \odot \varepsilon$, $\odot$ - element-wise product

spectrogram          • KL collapse problem

reference
encoder

$\downarrow$ FC  $\downarrow$ FC

$m$  $\sigma$

$z$

- the convergence speed of
  KL loss surpasses that of
  the reconstruction loss

· Several dimensions of latent $z$
  could independently control
  different style attributes

Total loss: $KL[q_\phi(z|x) \| p_\theta(z)] - \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x|z,t)] + l_{stop}$

# 3. Deep learning for Audio style transfer

· Main idea — learning content and style separately
  using <u>Gram Matrix</u> for style retrieval

and <u>convolutional NN's</u> for content

prepocessing — STFT (short time
Fourier transform)

therefore: spectogram

· Using different losses for style and content
(still a lot is unclear)

mean squared error
between Gram Matr. of image
& white noise image

mean squared err.
betw. source img. embed
& white noise embe.

Gram Matr. helps to numerically evaluate the correlation between features — the more they're similar the bigger the loss is.

So for <u>convolved</u> matr. the Gram matr. represents the "style", as it evaluates how these higher level features correlate { which repr. some "higher" level features, like, for images, edges, sharpness, blur, etc }

the total loss — $\sum$(style_loss, content_loss)

Important — same approach as with images can be used

4. Linear transformation Approaches to Many-to-One Voice Conversion

Voice conversion: 1) training
2) adaptation ( speaker selection, eigen voice techniques)
3) conversion

(1) Constrained Max. Likelyhood Linear Regression

$$\hat{x}_t = A x_t + b = W \zeta_t \quad \left\{ \begin{array}{l} W = [b, A] \\ \zeta_t = [1, x_t^\top]^\top \end{array} \right\}$$

~ constr. model-space transformation

$$P(x_t, Y_t \mid \lambda, W) = \sum_{m=1}^{M} \alpha_m \, \mathcal{N}\left( \begin{bmatrix} x_t \\ Y_t \end{bmatrix}; \hat{\mu}_m^{xx}, \hat{\Sigma}_m^{xx} \right)$$

where 

$$\hat{\mu}_m = \begin{bmatrix} A' \mu_m^x + b' \\ \mu_m^Y \end{bmatrix}$$

$$\hat{\Sigma}_m = \begin{bmatrix} A' \Sigma_m^{xx} A'^T & A' \Sigma_m^{xy} \\ \Sigma_m^{yx} A'^T & \Sigma_m^{yy} \end{bmatrix}$$

$\Rightarrow$ EM algorithm maximizes

$$\hat{W} = \underset{W}{\arg\max} \; \prod_{t=1}^{T} P(x_t^{(new)} \mid \lambda, W)$$

(2) Mean Linear Transformation

$$\hat{\mu}_m^x = A' \mu_m^x + b' = W' \zeta_m \qquad \begin{array}{l} W' = [b', A'] \\ \zeta_m = [1, \mu_m^{x\top}]^\top \end{array}$$

$$\sim P(x_t, Y_t \mid \lambda, W') = \sum_{m=1}^{M} \alpha_m \, \mathcal{N}\left( \begin{bmatrix} x_t \\ Y_t \end{bmatrix}; \hat{\mu}_m^{xy}, \Sigma_m^{xy} \right)$$

where $\hat{M}_m^{xy} = [M_m^{x\,T}, \; M_m^{y\,T}]^T$

And EM algorithm max.

$$\hat{W}' = \underset{W'}{argmax} \; \prod_{t=1}^{T} P(x_t^{(new)} \mid \lambda, W')$$

SAT = Speaker adaptive training

In SAT realisation with EM-algorithm both the params. of canonical GMM and a set of speaker-dependent linear transforms are optimised:

$$\{\hat{\lambda}, \hat{W}\} = \underset{\{\lambda, W\}}{argmax} \; \prod_{s=1}^{S} \prod_{t=1}^{T} P(X_t^{(s)}, Y_t \mid \lambda, W^{(s)})$$