

NEURAL VOICE CONVERSION AND IT'S RECENT ADVANCEMENTS

Dipjyoti Paul



University of Crete, Computer Science Dept., Greece,

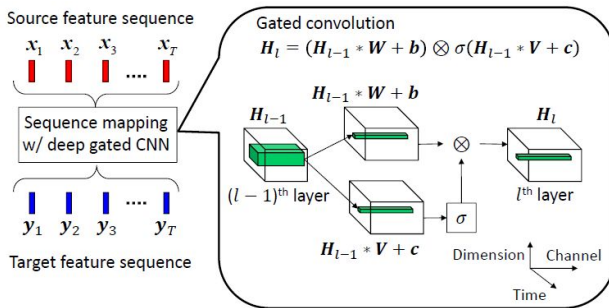
HY578: Digital Speech Signal Processing
Nov 2025



- 1 INTRODUCTION TO VC
- 2 NEURAL VOICE CONVERSION
 - Sequence based VC
 - VAE VC
 - GAN VC
 - zero-shot/few-shot VC
- 3 TTS TO VC
- 4 REFERENCES
- 5 THANK YOU

SEQUENCE BASED VC

- The proposed method uses Gated CNN to model long-term dependencies [Kaneko et. al. 2017].
- Training with fixed length of sequence.
- Inference with arbitrary length of sequence.



VARIATIONAL AUTOENCODER (VAE)-VC

- The core of VAE-VC is an encoder-decoder network.
- During training, given an observed (source or target) spectral frame \mathbf{x} , a speaker-independent encoder E_θ with parameter set θ encodes \mathbf{x} into a latent code:

$$\hat{\mathbf{z}} = E_{\theta}(\mathbf{x})$$

$$\hat{\mathbf{x}} = G_\phi(\hat{\mathbf{z}}, \hat{\mathbf{y}}) = G_\phi(E_\theta(\mathbf{x}), \hat{\mathbf{y}})$$

VARIATIONAL AUTOENCODER (VAE)-VC

- The core of VAE-VC is an encoder-decoder network.
- During training, given an observed (source or target) spectral frame \mathbf{x} , a speaker-independent encoder E_θ with parameter set θ encodes \mathbf{x} into a latent code:

$$\hat{\mathbf{z}} = E_\theta(\mathbf{x})$$

- The speaker code $\hat{\mathbf{y}}$ of the input frame is then concatenated with the latent code, and passed to a conditional decoder G_ϕ with parameter set ϕ to reconstruct the input.

$$\hat{\mathbf{x}} = G_\phi(\hat{\mathbf{z}}, \hat{\mathbf{y}}) = G_\phi(E_\theta(\mathbf{x}), \hat{\mathbf{y}})$$

VARIATIONAL AUTOENCODER (VAE)-VC

- The core of VAE-VC is an encoder-decoder network.
- During training, given an observed (source or target) spectral frame \mathbf{x} , a speaker-independent encoder E_θ with parameter set θ encodes \mathbf{x} into a latent code:

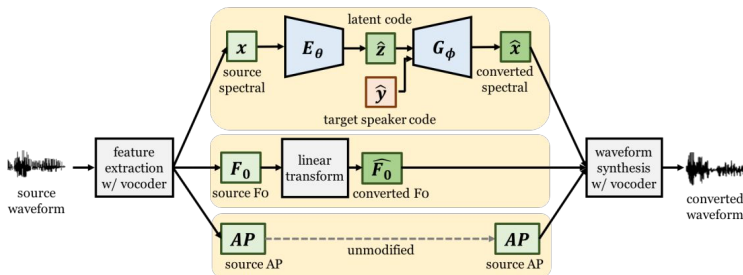
$$\hat{\mathbf{z}} = E_\theta(\mathbf{x})$$

- The speaker code $\hat{\mathbf{y}}$ of the input frame is then concatenated with the latent code, and passed to a conditional decoder G_ϕ with parameter set ϕ to reconstruct the input.

$$\hat{\mathbf{x}} = G_\phi(\hat{\mathbf{z}}, \hat{\mathbf{y}}) = G_\phi(E_\theta(\mathbf{x}), \hat{\mathbf{y}})$$

VAE-VC

- This framework is based on variational auto-encoder which exploits non-parallel corpora. [Hsu et. al. 2016]





VAE-VC

- The final approximated objective function of an individual frame:

$$\mathcal{L}(\theta, \phi; \mathbf{x}_n) = -D_{KL}(q_{\phi}(\mathbf{z}_n|\mathbf{x}_n)||p(\mathbf{z}_n)) + \mathbb{E}_{q_{\phi}(\mathbf{z}_n|\mathbf{x}_n)}[\log p_{\theta}(\mathbf{x}_n|\mathbf{z}_n, \hat{\mathbf{y}}_n)]$$

where, $q_{\phi}(\cdot)$ is the variational posterior.

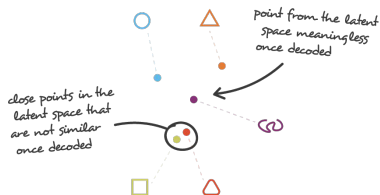
$p(\cdot)$ is the true posterior.

$D_{KL}(\cdot||\cdot)$ is the Kullback-Leibler divergence (KLD) of the approximate from the true posterior.

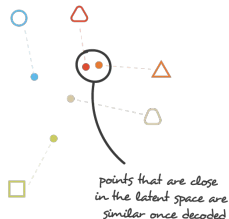
- Training is equivalent to iteratively finding the parameters that maximize the variational lower bound:

$$\{\theta^*, \phi^*\} = \operatorname{argmax}_{\theta, \phi} \mathcal{L}(\theta, \phi; \mathbf{X})$$

INTUITIONS ABOUT REGULARIZATION



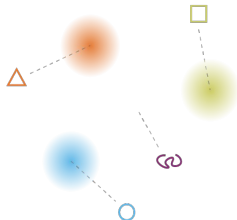
irregular latent space



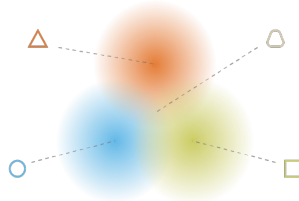
regular latent space



INTUITIONS ABOUT REGULARIZATION



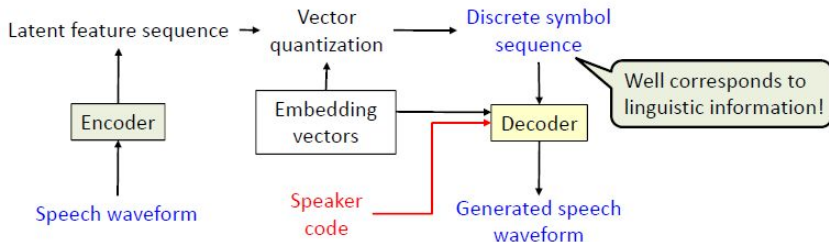
what can happen without regularisation



what we want to obtain with regularisation

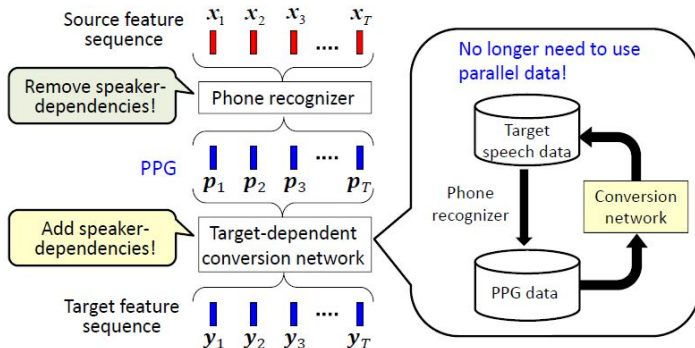
VECTOR QUANTIZATION VAE (VQ-VAE) VC

- Directly encodes speech waveform into a discrete symbol sequence capturing long-term dependencies. [Van den Oord et. al. 2017]
- The posterior and prior distributions are categorical
- The samples drawn from these distributions index an embedding table. These embeddings are then used as input into the decoder network.



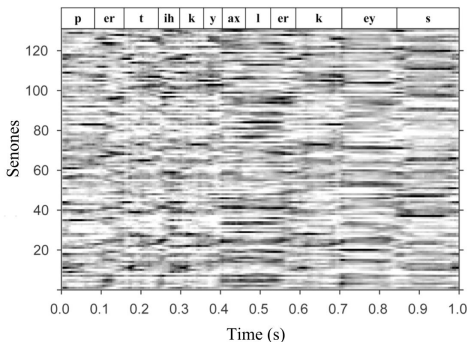
PHONEME POSTERIOGRAM VC

- Extract phoneme posteriorgram (PPG) as speaker independent contextual features. [Sun et. al. 2016]



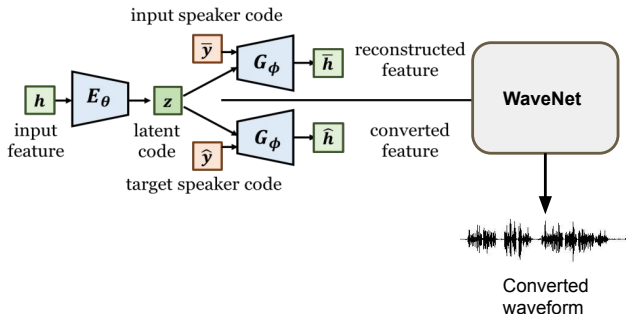
PHONEME POSTERIOGRAM VC

- PPG representation of the spoken phrase “particular case”.
- The horizontal axis (time in seconds), the vertical (indices of phonetic classes). The number of senones is 131. Darker shade implies a higher posterior probability.



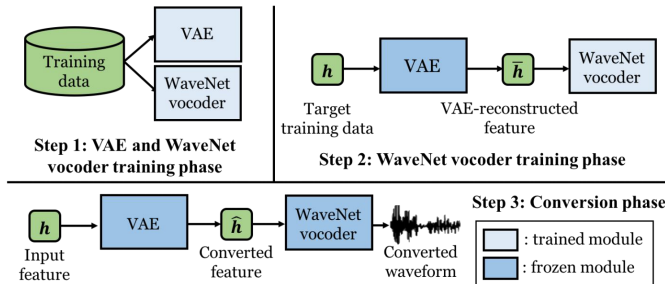
WAVENET VOCODER IN VAE-VC

- A general framework of WaveNet vocoder in voice conversion. [Huang et. al. 2019]



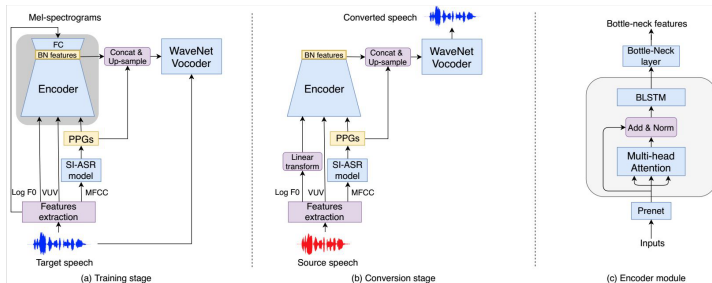
WAVENET VOCODER IN VAE-VC

• Training and Inference



JOINTLY TRAINED CONVERSION MODEL AND VOCODER

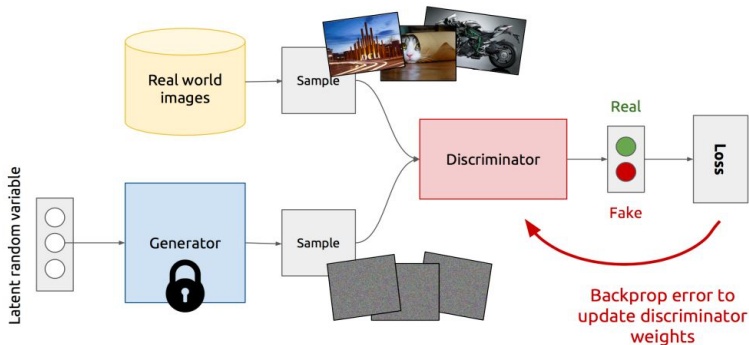
- The proposed model jointly trains a conversion model that maps phonetic posteriorgrams (PPGs) to Mel-spectrograms and a WaveNet vocoder. [Liu et. al. 2019]



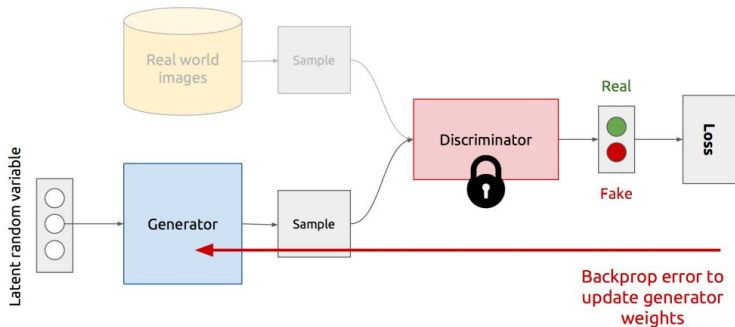
GAN FORMULATION

- An adversarial process which simultaneously trains two models: a generative model G that captures the data distribution, and a discriminative model D that estimates the probability that a sample came from the training data rather than G . [Goodfellow et. al. 2017]

DISCRIMINATOR (D) TRAINING



GENERATOR (G) TRAINING

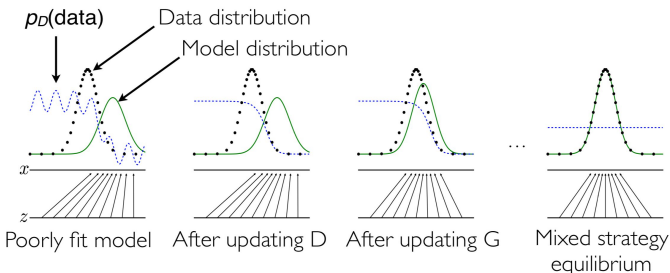


MATHEMATICAL NOTATIONS

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

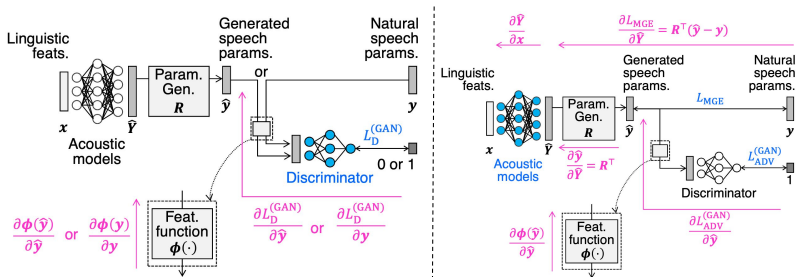
Value of \downarrow
 Expectation \swarrow
 prob. of $D(\text{real})$ \swarrow
 prob. of $D(\text{fake})$ \swarrow
 \min_G \max_D $V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$
 Minimize G \swarrow Maximize D \swarrow
 \mathbf{x} is sampled from real data \swarrow \mathbf{z} is sampled from $N(0, I)$ \swarrow fake

LEARNING GANs



GAN-BASED VC

- Training and gradients for updating the discriminator and Generator. ADV is adversarial loss and MGE refers to minimum generation loss. [Saito et. al. 2018]

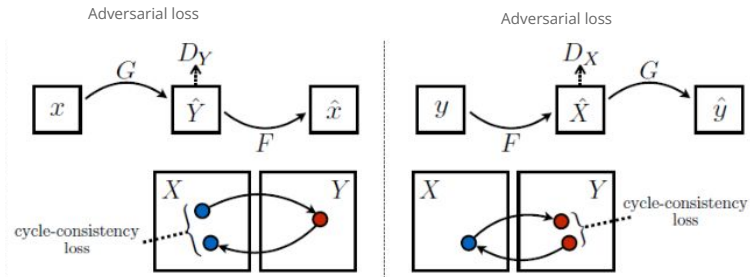


CYCLEGAN VOICE CONVERSION

- A non-parallel voice-conversion (VC) method that can learn a mapping from source to target speech without relying on parallel data. [Kaneko et. al. 2018]
- Forward and inverse mappings are simultaneously learned using an adversarial loss and cycle-consistency loss.
- Important losses are introduced:
 - Adversarial loss
 - Cycle-consistency loss
 - Identity-mapping loss

CYCLEGAN LOSSES

- The model learns mapping from source x to target y and vice versa.



CYCLEGAN LOSSES

- Two mapping function (Adversarial loss): G and F .

$$G : X \rightarrow Y \text{ and } F : Y \rightarrow X$$

- Cycle-consistency loss:
 - Forward: $x \rightarrow G(x) \rightarrow F(G(x)) \rightarrow \hat{x}$
 - Backward: $y \rightarrow F(y) \rightarrow G(F(y)) \rightarrow \hat{y}$
- Adversarial loss + Cycle-consistency loss

$$\mathcal{L}_{adv}(G_{X \rightarrow Y}, D_Y) + \mathcal{L}_{adv}(G_{Y \rightarrow X}, D_X) + \lambda_{cycle} \mathcal{L}_{cycle}(G_{X \rightarrow Y}, G_{Y \rightarrow X})$$

CYCLEGAN LOSSES

- Adversarial loss:

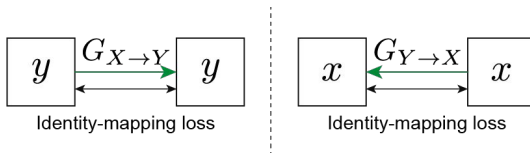
$$\begin{aligned}\mathcal{L}_{adv}(G_{X \rightarrow Y}, D_Y) &= \mathbb{E}_{y \sim P_{Data(y)}} [\log D_Y(y)] \\ &+ \mathbb{E}_{x \sim P_{Data(x)}} [\log(1 - D_Y(G_{X \rightarrow Y}(x)))]\end{aligned}$$

- Cycle-consistency loss:

$$\begin{aligned}\mathcal{L}_{cycle}(G_{X \rightarrow Y}, G_{Y \rightarrow X}) &= \mathbb{E}_{x \sim P_{Data(x)}} (\|G_{Y \rightarrow X}(G_{X \rightarrow Y}(x)) - x\|_1) \\ &+ \mathbb{E}_{y \sim P_{Data(y)}} (\|G_{X \rightarrow Y}(G_{Y \rightarrow X}(y)) - y\|_1)\end{aligned}$$

CYCLEGAN LOSSES

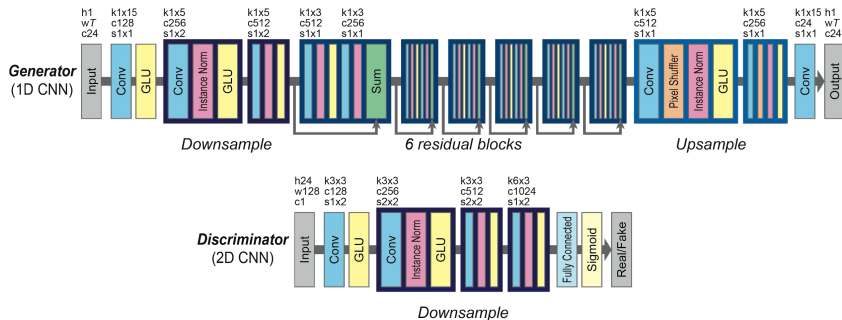
- Identity-mapping loss:
To encourage linguistic-information preservation, an identity-mapping loss is implemented.
- It encourages the generator to find the mapping that preserves composition between the input and output.



$$\mathcal{L}_{id}(G_{X \rightarrow Y}, G_{Y \rightarrow X}) = \mathbb{E}_{y \sim P_{Data(y)}} (\|G_{X \rightarrow Y}(y) - y\|_1) \\ + \mathbb{E}_{x \sim P_{Data(x)}} (\|G_{Y \rightarrow X}(x) - x\|_1)$$

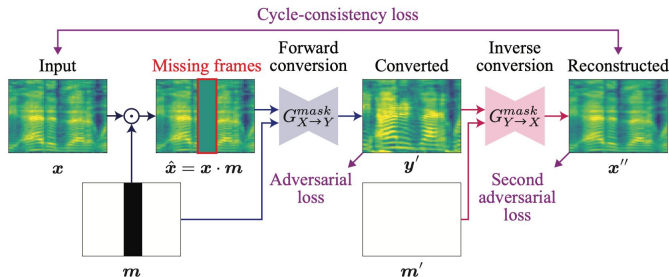
CYCLEGAN ARCHITECTURE

- The model architecture.



MASKED CYCLEGAN VC

- It is trained using a novel auxiliary task called filling in frames.
- A temporal mask to the input mel-spectrogram and encourage the converter to fill in missing frames.
- Allows the converter to learn time-frequency structures in a self-supervised manner. [Kaneko et. al. 2021]



STARGAN LOSSES

- The full objectives of StarGAN-VC to be minimized with respect to G , D and C are

- Generator loss:

$$\mathcal{L}_{adv}(G) + \lambda_{cls}\mathcal{L}_{cls}(G) + \lambda_{cyc}\mathcal{L}_{cyc}(G) + \lambda_{id}\mathcal{L}_{id}(G)$$

- Discriminator loss

$$\mathcal{L}_{adv}(D)$$

- Classifier loss

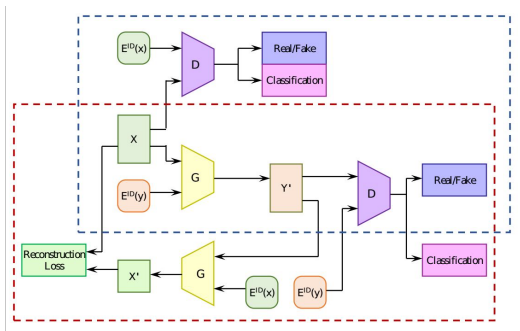
$$\mathcal{L}_{cls}(C)$$

ZERO-SHOT/FEW-SHOT VC

- The target speaker is unseen (zero-shot) during training or a very limited set of samples are available (few-shot).
- An universal embedding vector is used to represent speaker id.
- The idea is to represent any arbitrary unseen speaker ID with an embedding vector.
- Such embedding vector represents unseen speaker's timbre and would be a weighted combination of the timbres the speakers seen in the dataset. [Wang et. al. 2020]

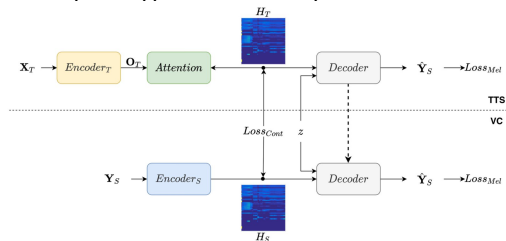
ZERO-SHOT STARGAN VC

- Illustration of our proposed StarGAN framework.
 - $E_{ID}(\cdot)$ represents the embedding ID of a speaker.
 - X' and Y' refer to the features reconstructed through G with embedded ID $E_{ID}(X)$ and $E_{ID}(Y)$ respectively.



TEXT-TO-SPEECH SYNTHESIS TO VOICE CONVERSION

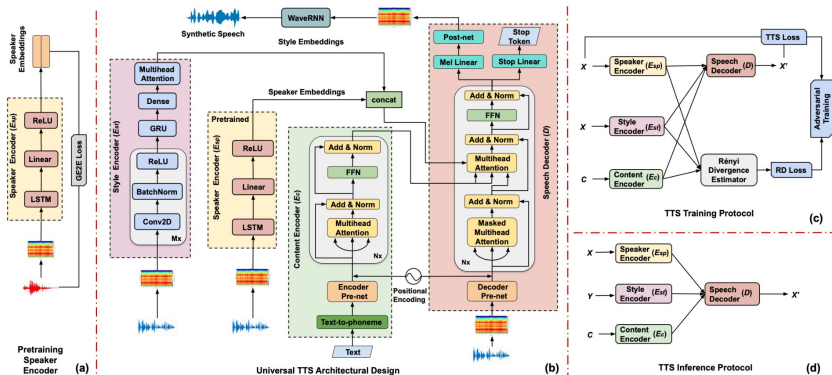
- VC framework by learning from a TTS synthesis system.
- The decoder is condition on a speaker embedding, becoming any-to-any VC.
- X_T denotes the input text, Y_S and \hat{Y}_S are target melspecs and the melspecs generated by the pipelines; O_T denotes the text encoding, H_T denotes the context vectors from TTS pipeline, H_S denotes the context vectors equivalents from the VC pipeline. [Zhang et. al. 2021]



MULTI-SPEAKER MULTI-STYLE TTS

- To generalize the models with multiple speakers and multiple styles using just the reconstruction loss, performance unfortunately deteriorates.
- Leaking content information into the style embeddings (“content leakage”) and leaking speaker information into style embeddings (“style leakage”).
- Minimizes the Reyni Diverence Disentagle Representation (RDDR) between the joint distribution and the product of marginals for the content-style and style-speaker pairs.
- Training is conducted on non-parallel data and generates voices in an unsupervised manner, i.e., neither style annotation nor speaker label are required. [Paul et. al. 2021]

MULTI-SPEAKER MULTI-STYLE TTS



MULTI-SPEAKER MULTI-STYLE TTS

Algorithm: Pseudo-code for proposed RDDR training

Input: Speech and text pairs $\langle \mathbf{x}_i, \mathbf{c}_i \rangle$.

Pre-training: Optimize E_c, D on LJSpeech using

$$\min_{E_c, E_{st}, D} \sum_i \| D(E_c(\mathbf{c}_i), E_{st}(\mathbf{x}_i), E_{sp}(\mathbf{x}_i)) - \mathbf{x}_i \|_1$$

$E_{sp} \leftarrow$ GE2E training

$E_{st}, T_\theta, T'_{\theta'} \leftarrow$ initialization with random weights

while $E_{st}, D, T_\theta, T'_{\theta'}$ not converged **do**

 Sample mini-batch from $\langle \mathbf{x}_i, \mathbf{c}_i \rangle$; $i = \{1, 2, \dots, b\}$

$\{\mathbf{p}_i\} \leftarrow \{E_c(\mathbf{c}_i) | i = 1, 2, \dots, b\}$

$\{\mathbf{q}_i\} \leftarrow \{E_{st}(\mathbf{x}_i) | i = 1, 2, \dots, b\}$

$\{\mathbf{r}_i\} \leftarrow \{E_{sp}(\mathbf{x}_i) | i = 1, 2, \dots, b\}$

$\{\hat{\mathbf{p}}_i\}, \{\tilde{\mathbf{r}}_i\} \leftarrow$ random permutation of $\{\mathbf{p}_i\}, \{\mathbf{r}_i\}$

$$\mathcal{L}_{RD^1} = \sum_k \left[-\frac{1}{\beta_k} \log \frac{1}{b} \sum_{i=1}^b e^{-\beta_k T_\theta(\mathbf{p}_i, \mathbf{q}_i)} \right.$$

$$\left. - \frac{1}{\gamma_k} \log \frac{1}{b} \sum_{i=1}^b e^{\gamma_k T_\theta(\hat{\mathbf{p}}_i, \mathbf{q}_i)} \right]$$

$$\mathcal{L}_{RD^2} = \sum_k \left[-\frac{1}{\beta_k} \log \frac{1}{b} \sum_{i=1}^b e^{-\beta_k T'_{\theta'}(\mathbf{r}_i, \mathbf{q}_i)} \right.$$

$$\left. - \frac{1}{\gamma_k} \log \frac{1}{b} \sum_{i=1}^b e^{\gamma_k T'_{\theta'}(\tilde{\mathbf{r}}_i, \mathbf{q}_i)} \right]$$

The overall objective function:

$$\mathcal{L} = \frac{1}{b} \sum_{i=1}^b \| D(\mathbf{p}_i, \mathbf{q}_i, \mathbf{r}_i) - \mathbf{x}_i \|_1$$

$$+ \lambda \max(0, \mathcal{L}_{RD^1}) + \lambda \max(0, \mathcal{L}_{RD^2})$$

$$D = D - \epsilon \nabla_D \mathcal{L}; \quad E_{st} = E_{st} - \epsilon \nabla_{E_{st}} \mathcal{L}$$

$$T_\theta = T_\theta + \epsilon \nabla_D \mathcal{L}_{RD^1}; \quad T'_{\theta'} = T'_{\theta'} + \epsilon \nabla_D \mathcal{L}_{RD^2}$$

end

REFERENCES

- 9 Saito, Yuki, Shinnosuke Takamichi, and Hiroshi Saruwatari. "Statistical parametric speech synthesis incorporating generative adversarial networks." *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 26, no. 1 (2017): 84-96.
- 10 Kaneko, Takuhiro, and Hirokazu Kameoka. "Parallel-data-free voice conversion using cycle-consistent adversarial networks." *arXiv preprint arXiv:1711.11293* (2017).
- 11 Kaneko, Takuhiro, Hirokazu Kameoka, Kou Tanaka, and Nobukatsu Hojo. "Maskcyclegan-vc: Learning non-parallel voice conversion with filling in frames." In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 5919-5923. IEEE, 2021.
- 12 Kameoka, Hirokazu, Takuhiro Kaneko, Kou Tanaka, and Nobukatsu Hojo. "Stargan-vc: Non-parallel many-to-many voice conversion using star generative adversarial networks." In *2018 IEEE Spoken Language Technology Workshop (SLT)*, pp. 266-273. IEEE, 2018.
- 13 Wang, Ruobai, Yu Ding, Lincheng Li, and Changjie Fan. "One-shot voice conversion using star-GAN." In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 7729-7733. IEEE, 2020.
- 14 Zhang, Mingyang, Yi Zhou, Li Zhao, and Haizhou Li. "Transfer learning from speech synthesis to voice conversion with non-parallel training data." *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 29 (2021): 1290-1302.
- 15 Paul, Dipjyoti, Sankar Mukherjee, Yannis Pantazis, and Yannis Stylianou. "A Universal Multi-Speaker Multi-Style Text-to-Speech via Disentangled Representation Learning Based on Rényi Divergence Minimization." In *Interspeech*, pp. 3625-3629. 2021.
- 16 Tomoki Toda: Advanced Voice Conversion, Speech Processing Courses in Crete (SPCC2018).

THANK YOU
for your attention