

---

---

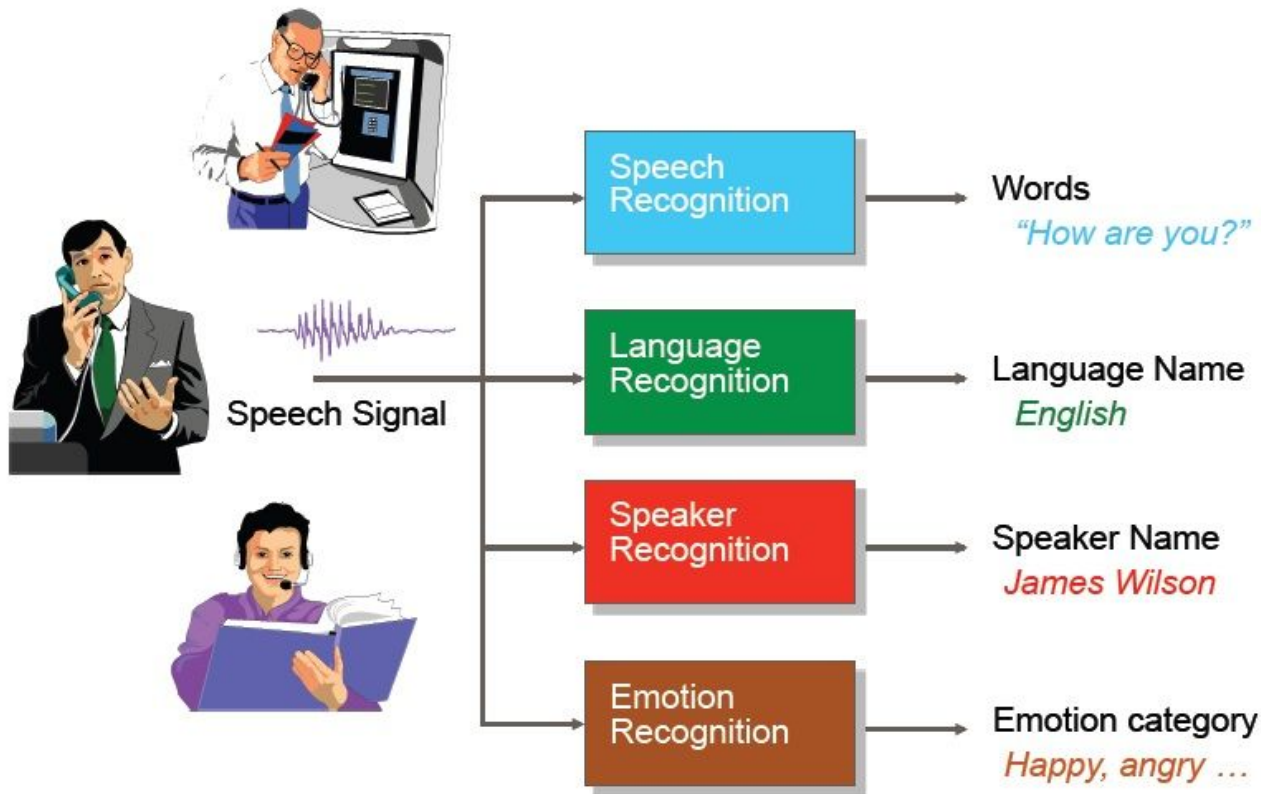
# Voice Conversion

— Dipjyoti Paul —  
University of Crete, Greece

---

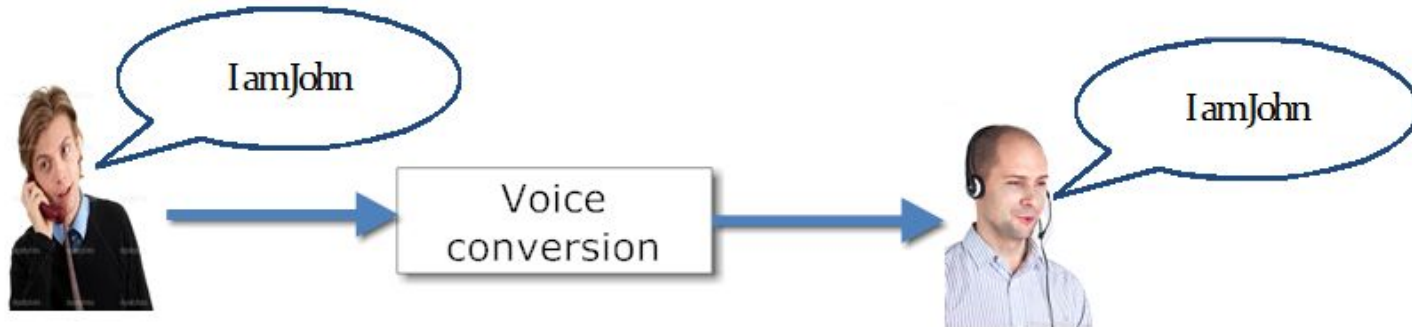
---

# Information in Speech



# Voice Conversion (VC)

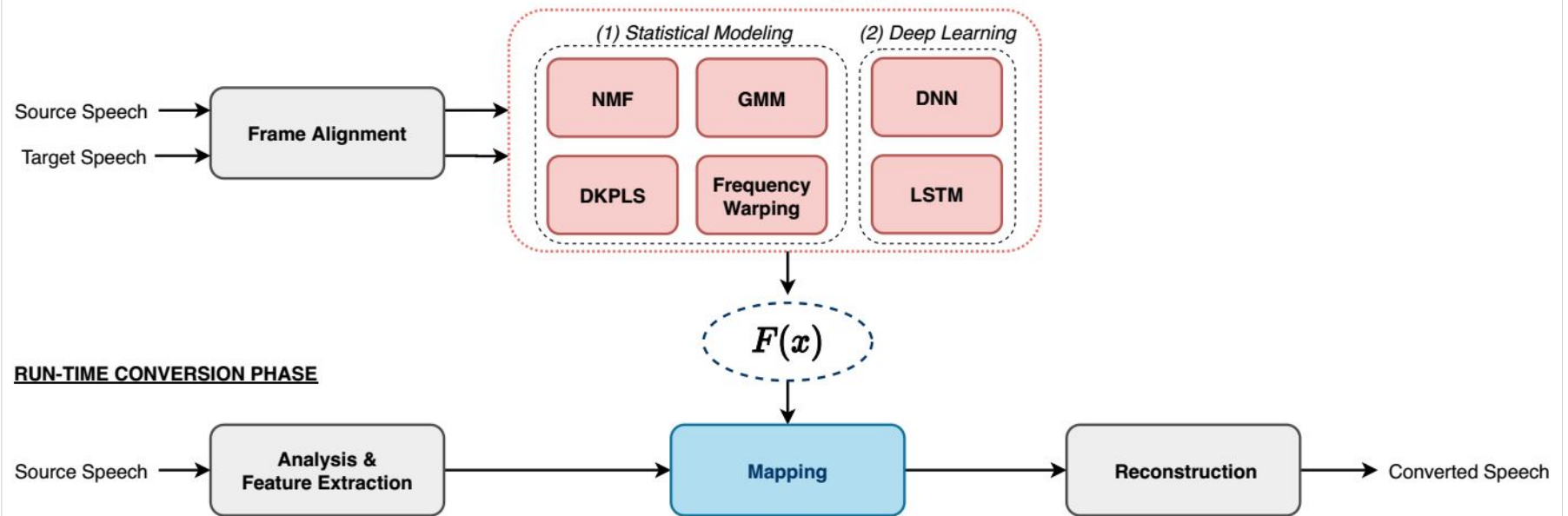
- Technique to convert the utterance of a source speaker to create the perception as if spoken by a specified target speaker.
- Only transform the speaker timbre (para-linguistic information) and keep the linguistic message in the utterance unchanged.



Source speaker's voice

Target speaker's voice

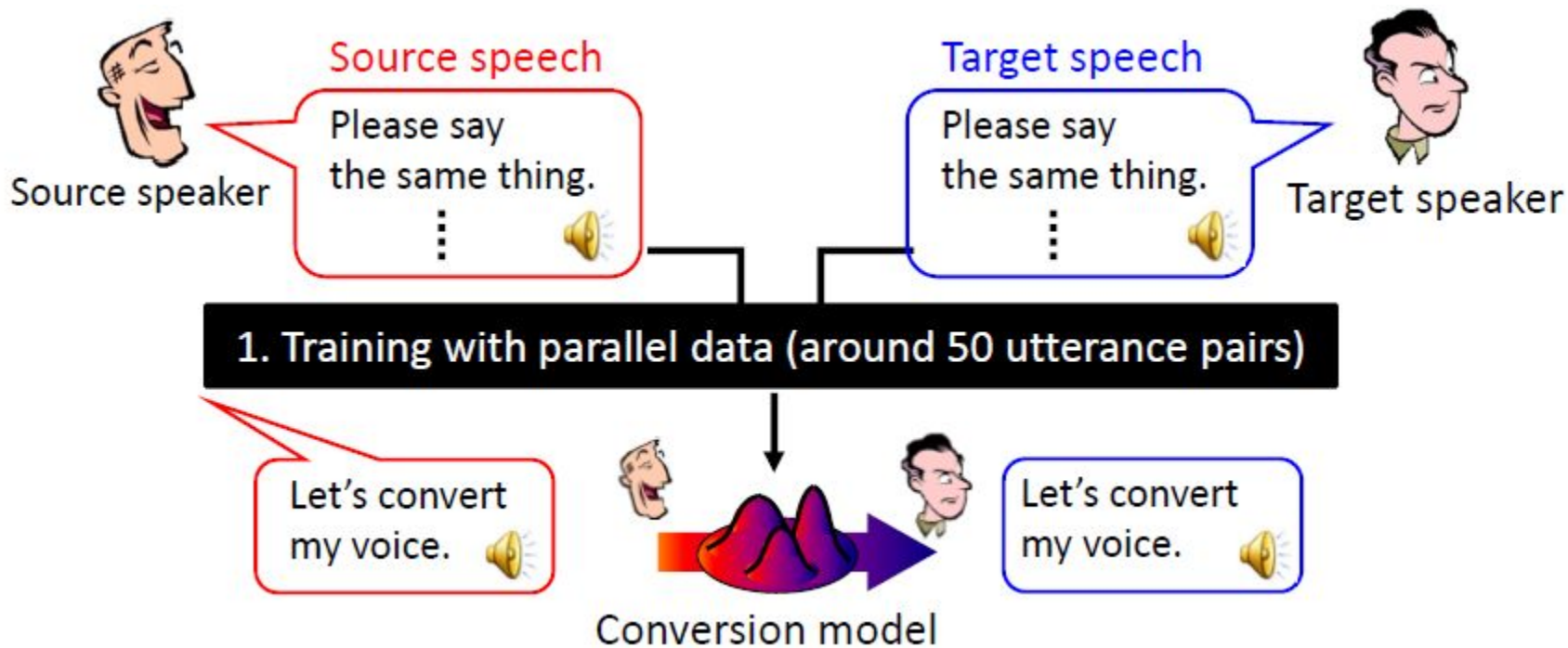
# Voice Conversion



# Applications

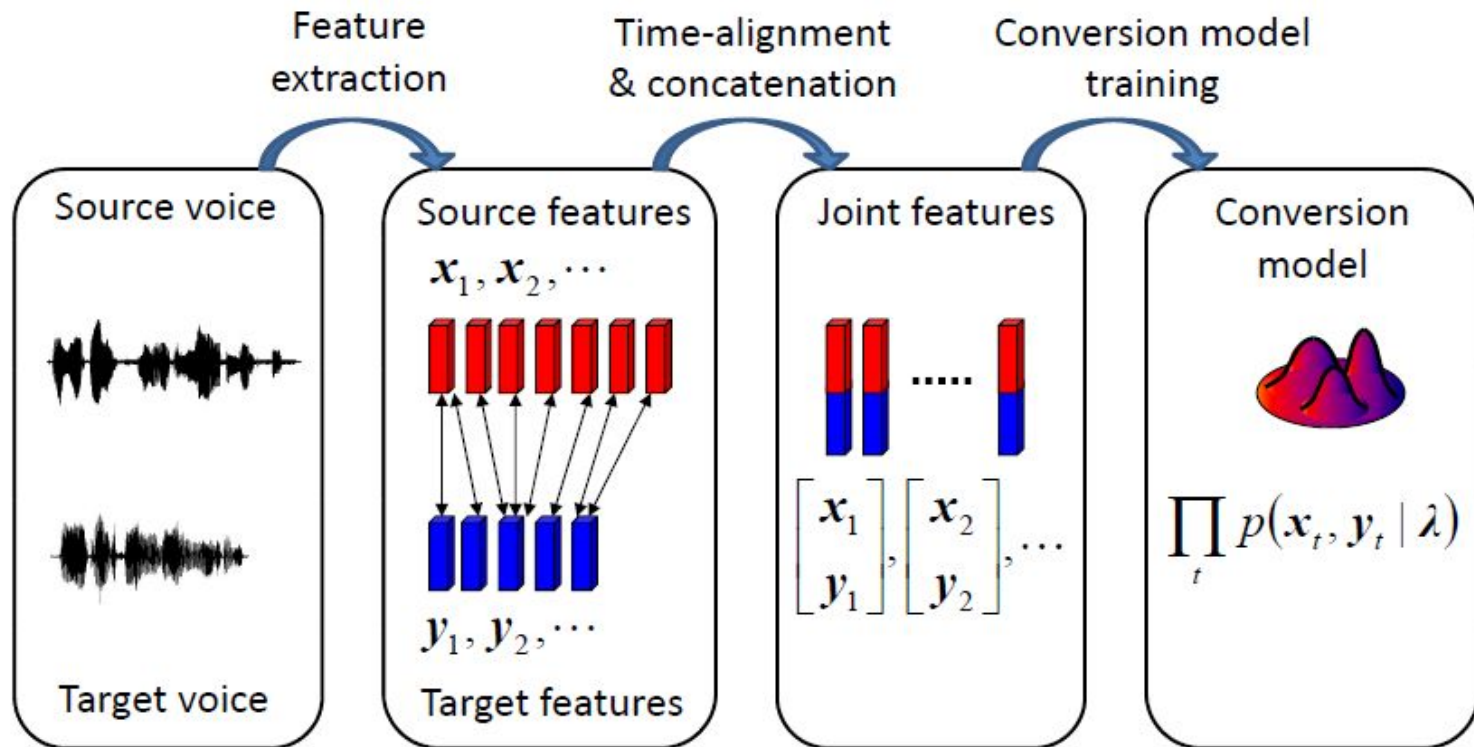
- Text-To-Speech (TTS) customization
- Film dubbing
- Design of speaking aids
- Education etc

# Statistical VC

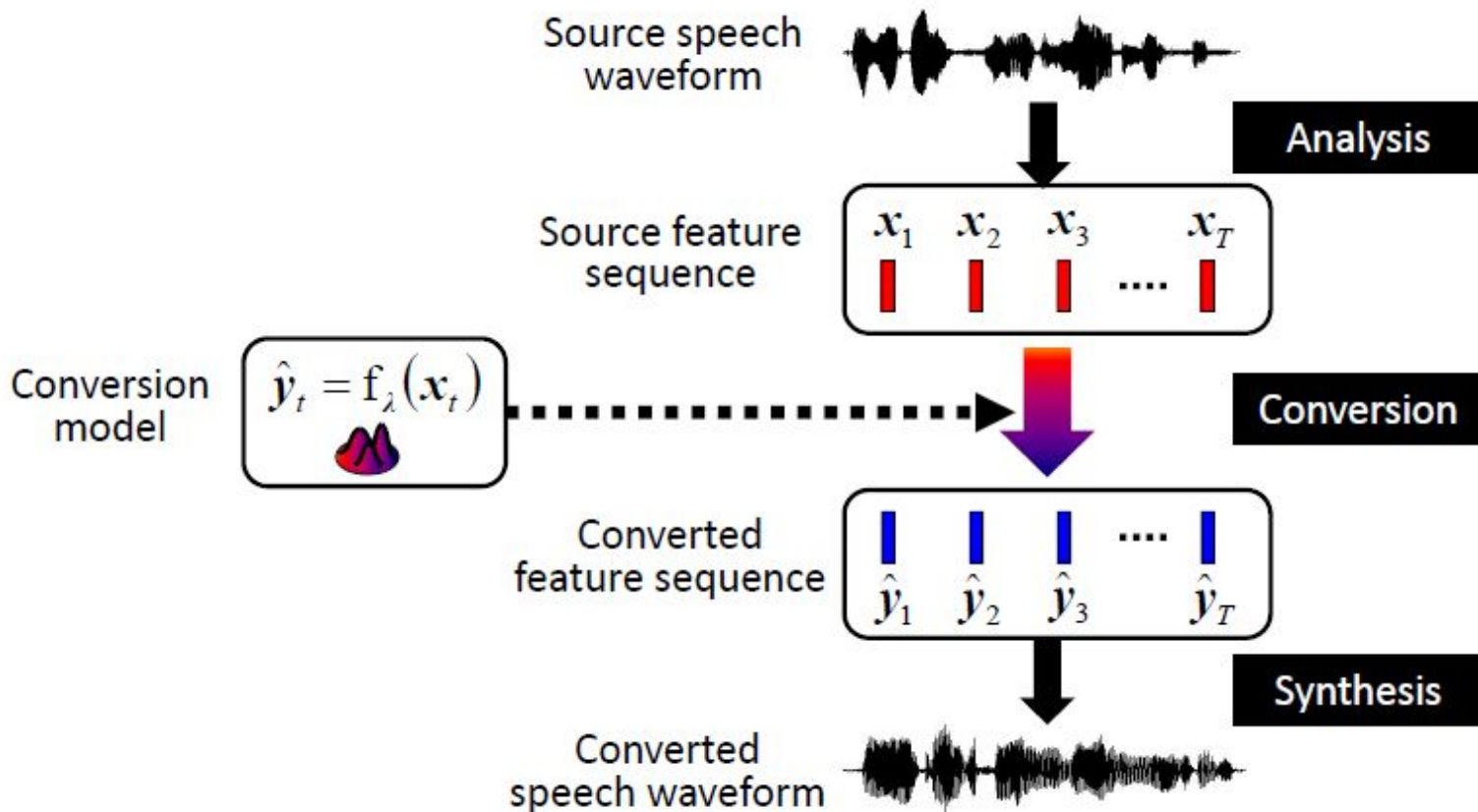


**2. Conversion of any utterance while keeping linguistic contents unchanged**

# VC Training

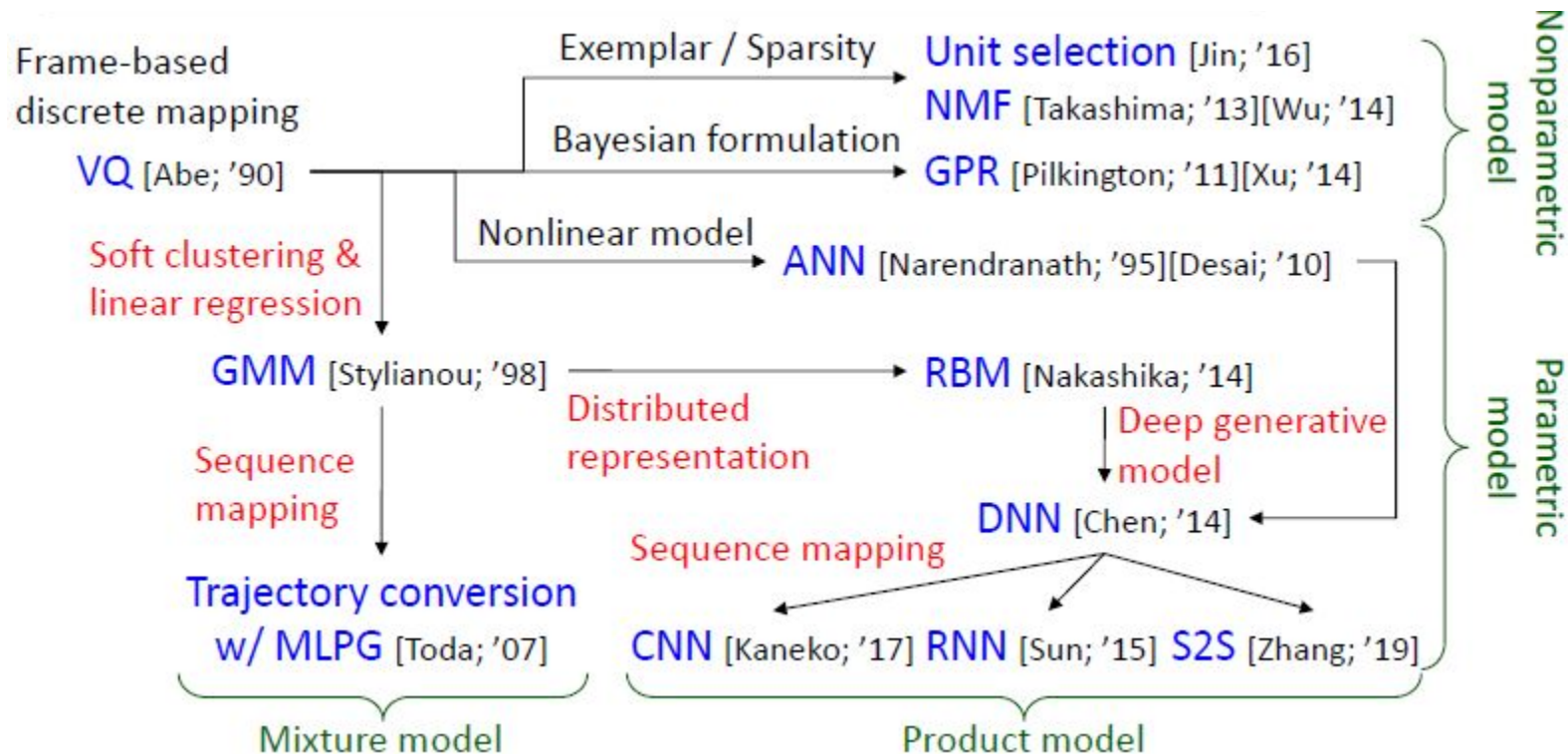


# VC Conversion





# Timeline of VC Research

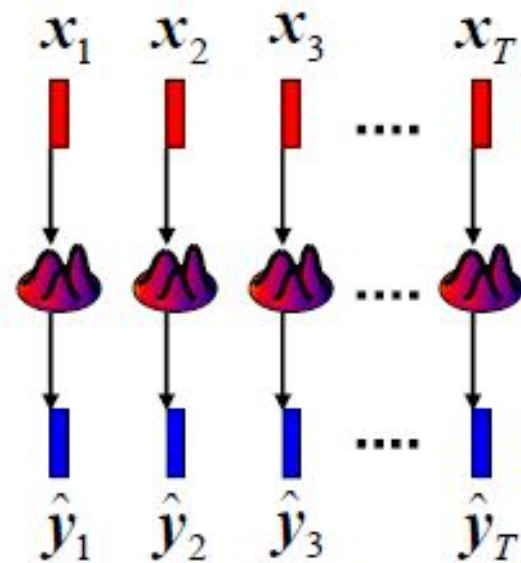


# Frame-based VC

- Source feature:  $\mathbf{x}$
- Target feature:  $\mathbf{y}$
- Converted feature:  $\hat{\mathbf{y}}$

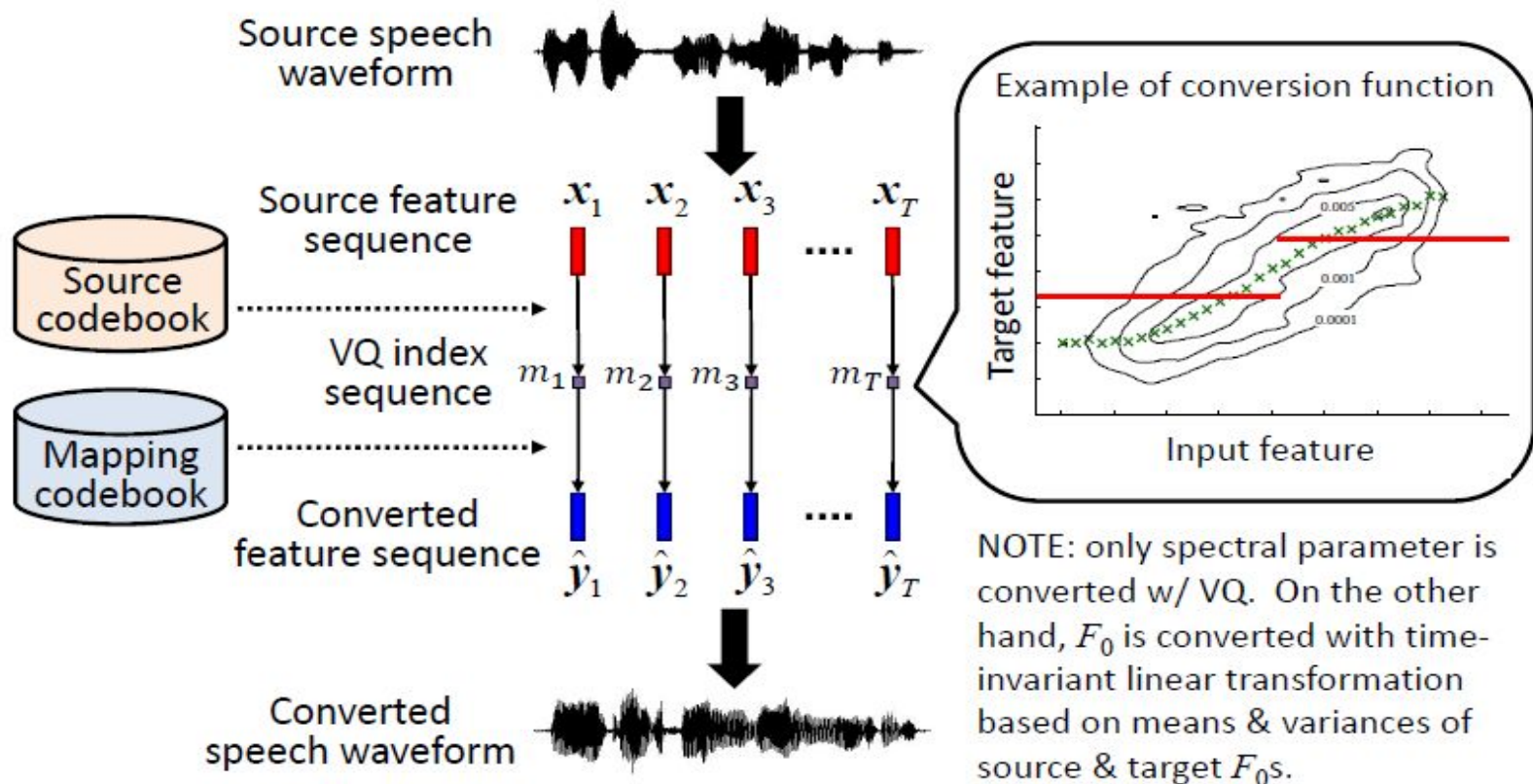
Frame-based conversion function

$$\hat{\mathbf{y}}_t = \mathbf{f}_\lambda(\mathbf{x}_t)$$



# Vector Quantization-based VC

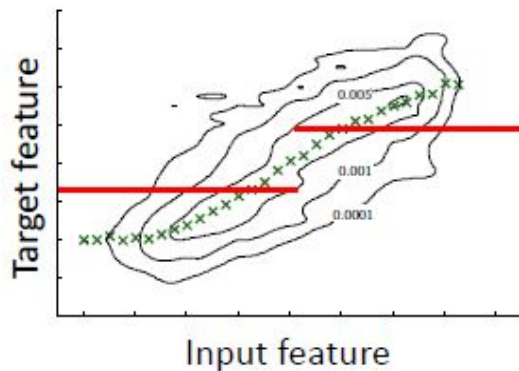
[Abe et. al. 1990]



# Discontinuous to Continuous Conversion

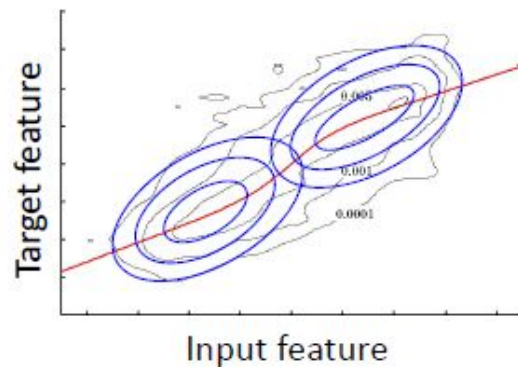
## VQ-based conversion

- Discrete function w/ hard clustering
- Ignore feature correlation w/ discrete mapping




## GMM-based conversion

- Continuous function w/ soft clustering
- Directly model feature correlation w/ linear regression



# GMM based Conversion

Joint feature vector:  $\mathbf{z}_t = \begin{bmatrix} \mathbf{x}_t \\ \mathbf{y}_t \end{bmatrix}$  

GMM:

$$\text{Joint } p.d.f.: P(\mathbf{x}_t, \mathbf{y}_t | \lambda) = \sum_{m=1}^M \alpha_m \mathcal{N}(\mathbf{z}_t; \mu_m^{(z)}, \Sigma_m^{(zz)})$$

Maximum likelihood training

$$\hat{\lambda} = \arg \max \prod_t P(\mathbf{x}_t, \mathbf{y}_t | \lambda)$$

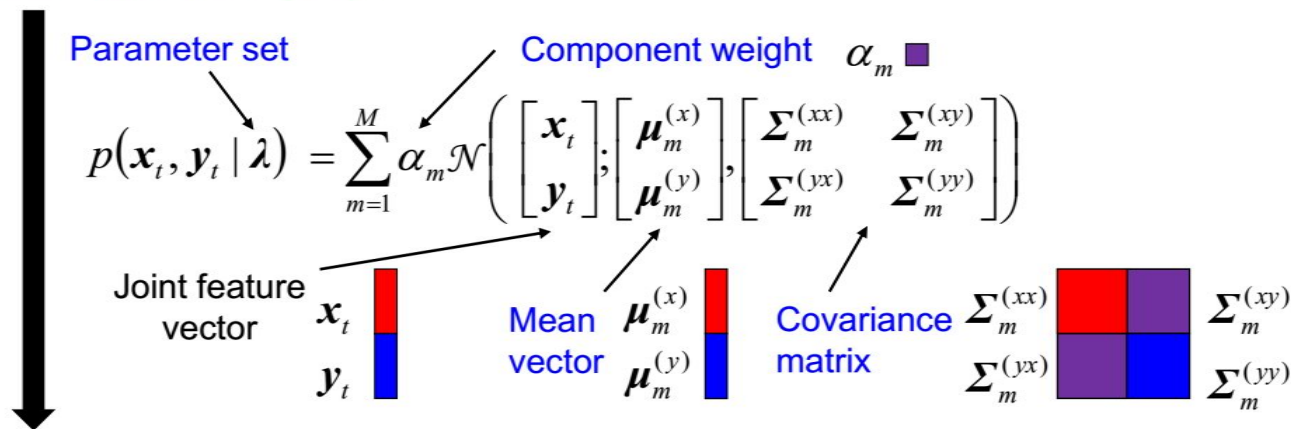
Updated model parameters

Likelihood for all joint vectors

# GMM based Conversion

[Stylianou et. al. 1998]

Training of joint *p.d.f.* (modeled by a GMM) [Kain; '98]



Conversion w/ conditional *p.d.f.* (also modeled by a GMM)

$$p(\mathbf{y}_t | \mathbf{x}_t, \lambda) = \frac{p(\mathbf{x}_t, \mathbf{y}_t | \lambda)}{\int p(\mathbf{x}_t, \mathbf{y}_t | \lambda) d\mathbf{y}_t} = \sum_{m=1}^M p(m | \mathbf{x}_t, \lambda) \mathcal{N}(\mathbf{y}_t; \boldsymbol{\mu}_{m,t}^{(y|x)}, \boldsymbol{\Sigma}_m^{(y|x)})$$

$$\text{MMSE estimate: } \hat{\mathbf{y}}_t = \int \mathbf{y}_t p(\mathbf{y}_t | \mathbf{x}_t, \lambda) d\mathbf{y}_t = \sum_{m=1}^M p(m | \mathbf{x}_t, \lambda) \boldsymbol{\mu}_{m,t}^{(y|x)}$$



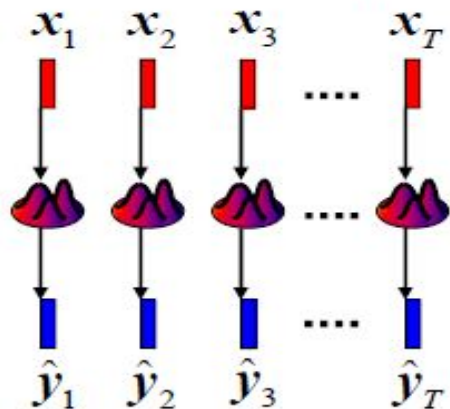
# Sequence-based VC

[Toda et. al. 2007]

## Frame-based conversion

$$\hat{y}_t = f_\lambda(x_t)$$

Source feature sequence

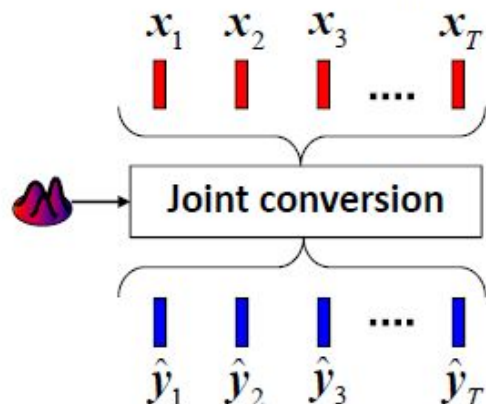


Converted feature sequence

## Sequence-based conversion

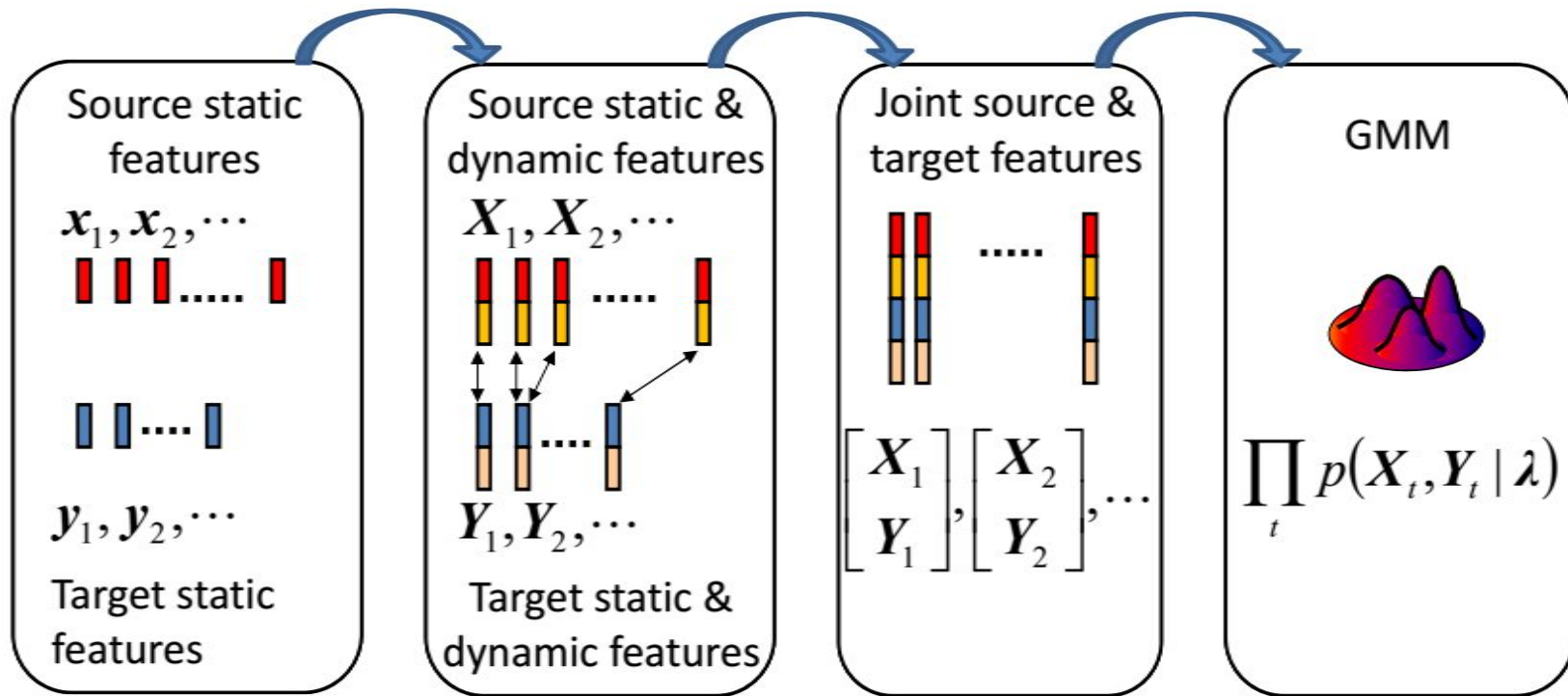
$$\{\hat{y}_1, \hat{y}_2, \dots, \hat{y}_T\} = f_\lambda(x_1, x_2, \dots, x_T)$$

Source feature sequence



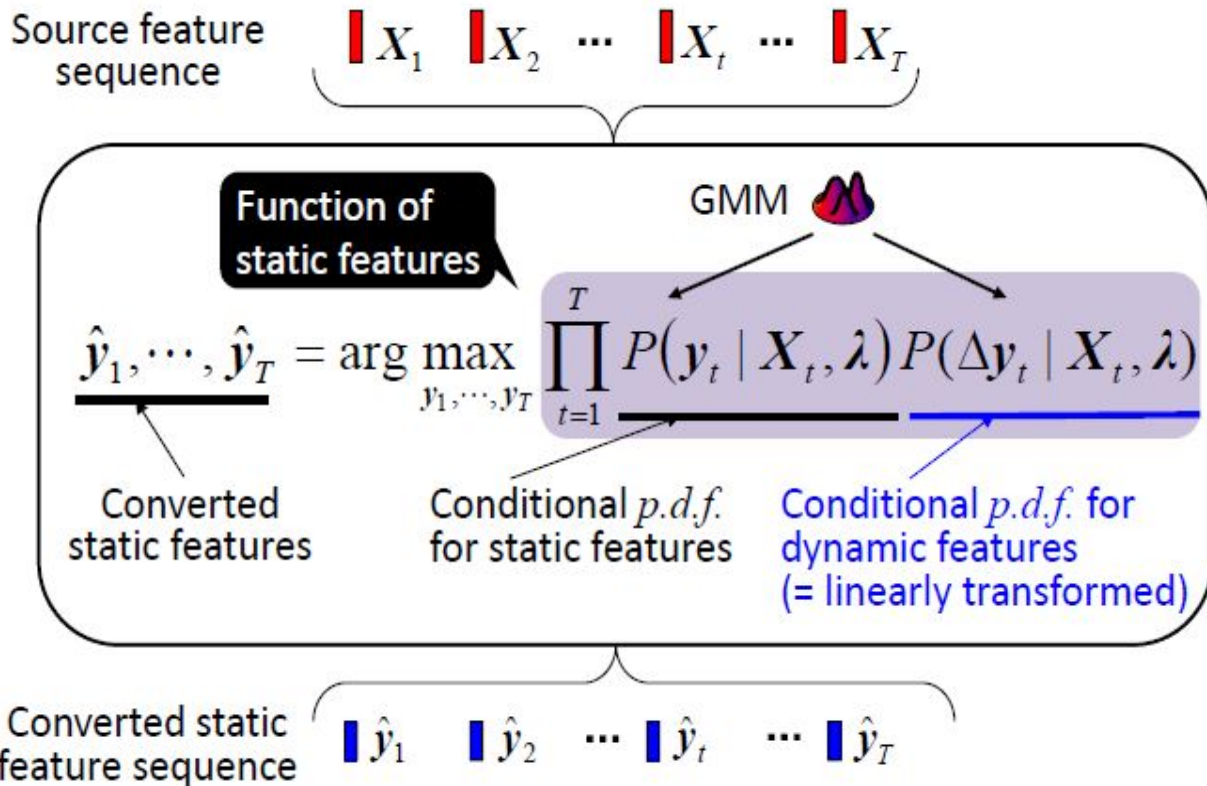
Converted feature sequence

# Sequence-based VC





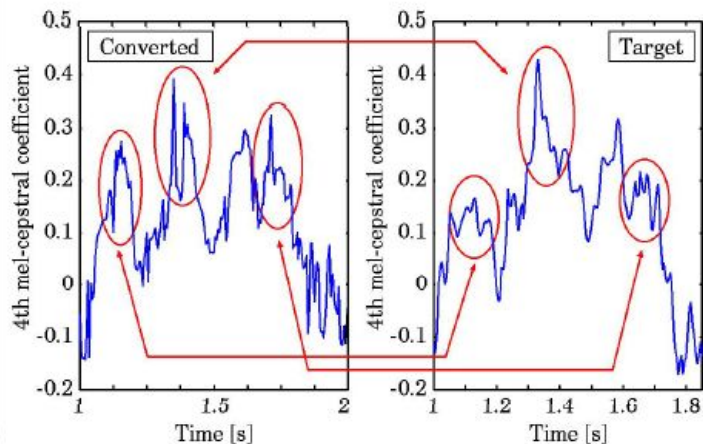
# Sequence-based VC



# Limitations of JD-GMM

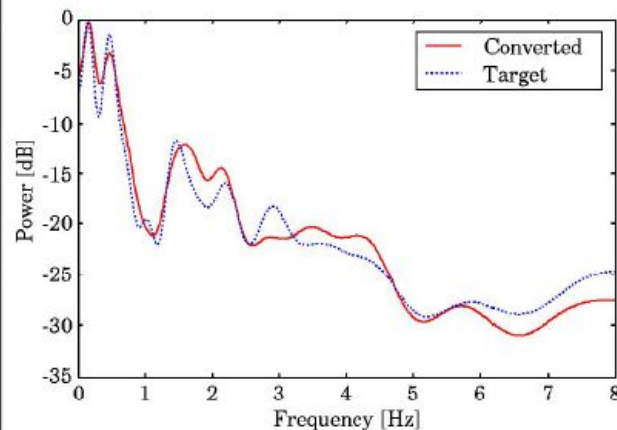
## Discontinuous transitions

- Ignoring inter-frame correlation...



## Over-smoothing effect

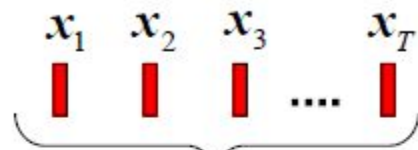
- Missing some characteristics not well modeled by GMM...



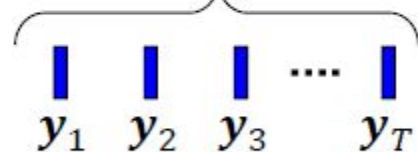
# VC based on Deep Neural Networks

# Sequence-based VC

Source feature sequence



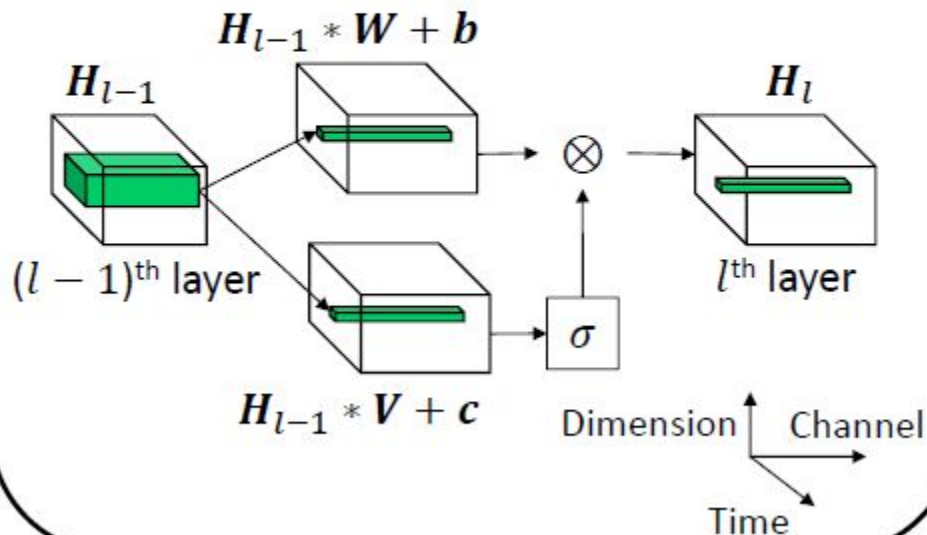
Sequence mapping  
w/ deep gated CNN



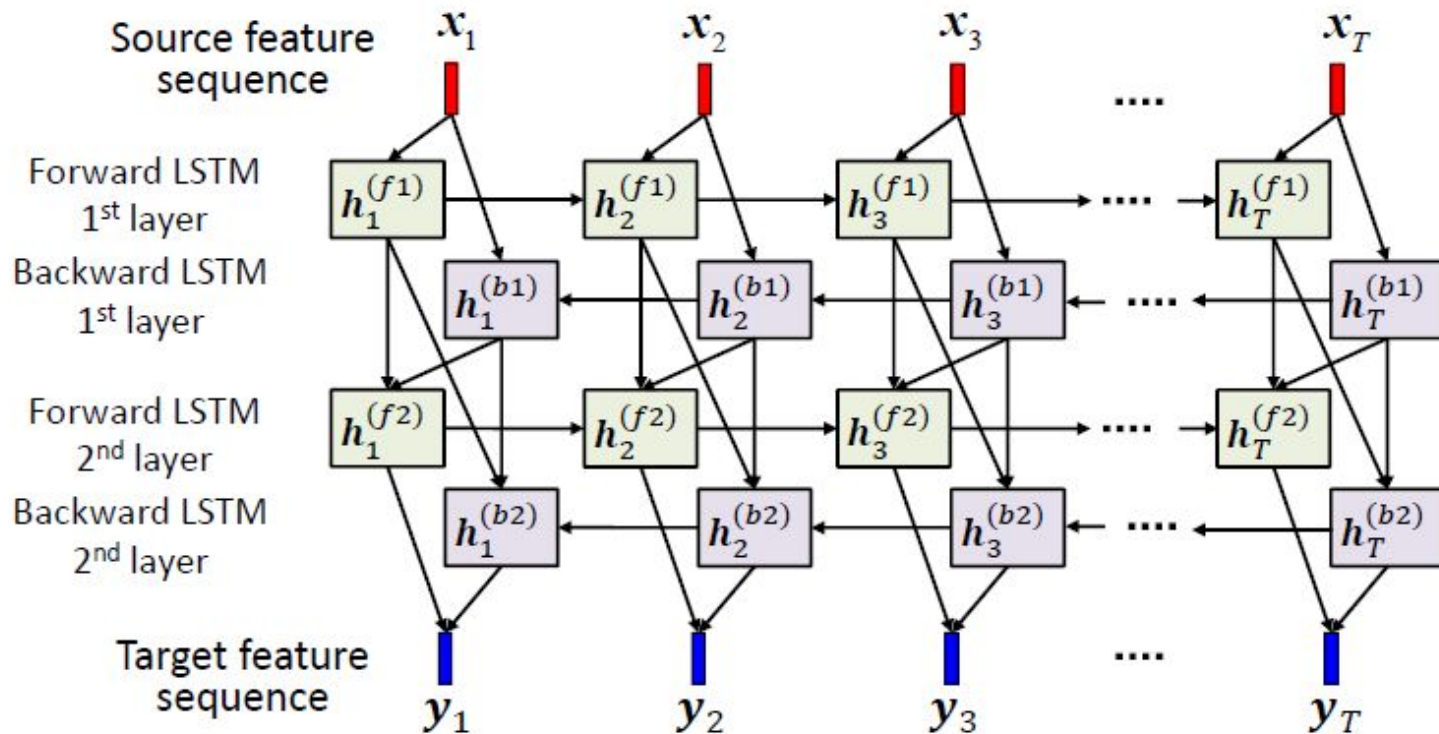
Target feature sequence

Gated convolution

$$H_l = (H_{l-1} * W + b) \otimes \sigma(H_{l-1} * V + c)$$



# Sequence-based VC

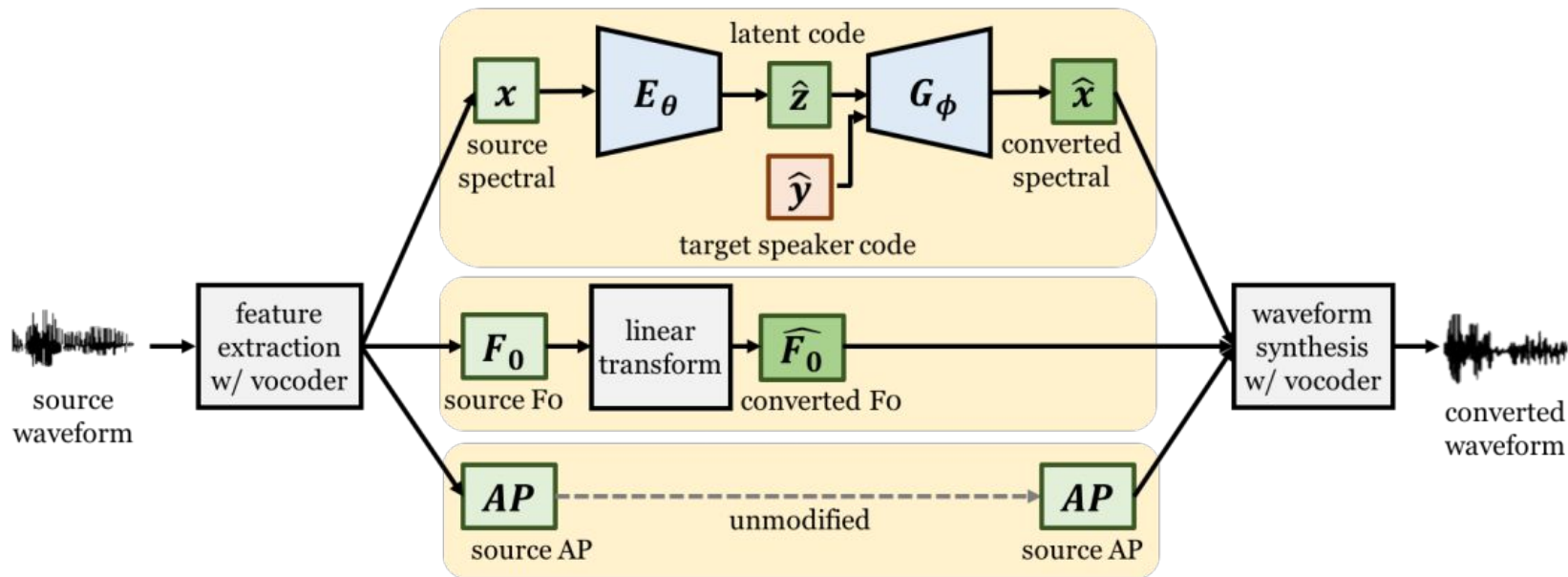


# Variational Autoencoder (VAE)-VC

- 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  $\mathbf{E}_\theta$  with parameter set  $\theta$  encodes  $\mathbf{x}$  into a latent code:  $\bar{\mathbf{z}} = E_\theta(\mathbf{x})$ .
- The speaker code  $\mathbf{y}$  of the input frame is then concatenated with the latent code, and passed to a conditional decoder  $\mathbf{G}_\phi$  with parameter set  $\phi$  to reconstruct the input.

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

# VAE-VC



# VAE-VC

- The model parameters can be obtained by maximizing the variational lower bound:

$$\begin{aligned}\mathcal{L}_{vae}(\theta, \phi; \mathbf{x}, \mathbf{y}) &= \mathcal{L}_{recon}(\mathbf{x}, \mathbf{y}) + \mathcal{L}_{lat}(\mathbf{x}), \\ \mathcal{L}_{recon}(\mathbf{x}, \mathbf{y}) &= \mathbb{E}_{\mathbf{z} \sim q_{\theta}(\bar{\mathbf{z}}|\mathbf{x})} [\log p_{\phi}(\bar{\mathbf{x}}|\mathbf{z}, \mathbf{y})], \\ \mathcal{L}_{lat}(\mathbf{x}) &= -D_{KL}(q_{\theta}(\bar{\mathbf{z}}|\mathbf{x}) \| p(\mathbf{z})),\end{aligned}$$

$q_{\theta}(\bar{\mathbf{z}}|\mathbf{x})$ : approximate posterior.

$p_{\phi}(\bar{\mathbf{x}}|\mathbf{z}, \mathbf{y})$ : data likelihood.

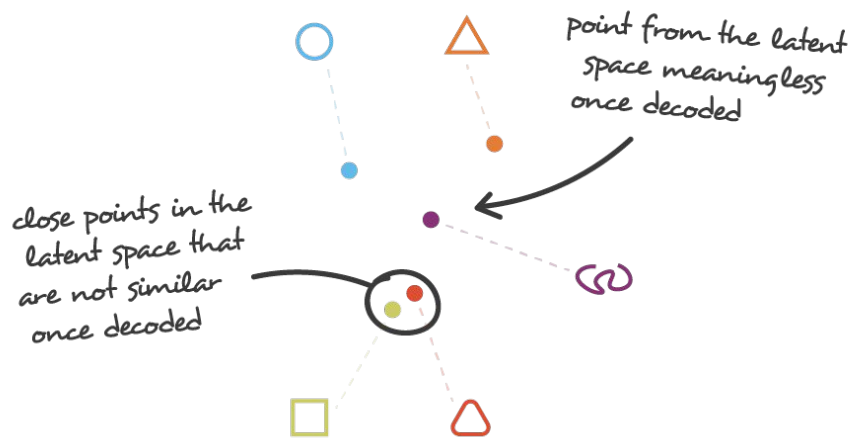
$p(\mathbf{z})$ : prior distribution of the latent space.

- Conversion phase:

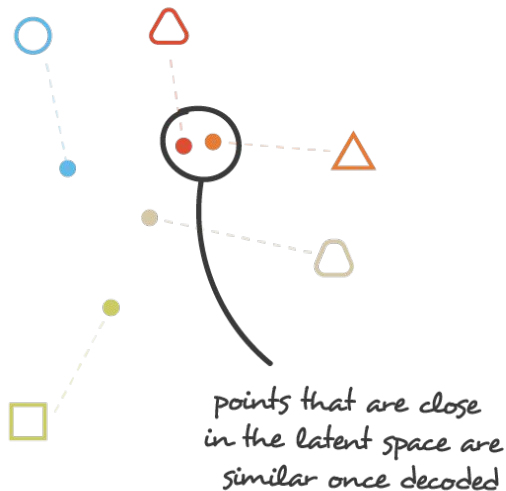
$$\hat{\mathbf{x}} = f(\mathbf{x}, \hat{\mathbf{y}}) = G_{\phi}(\hat{\mathbf{z}}, \hat{\mathbf{y}}) = G_{\phi}(E_{\theta}(\mathbf{x}), \hat{\mathbf{y}})$$



# Intuitions about Regularization

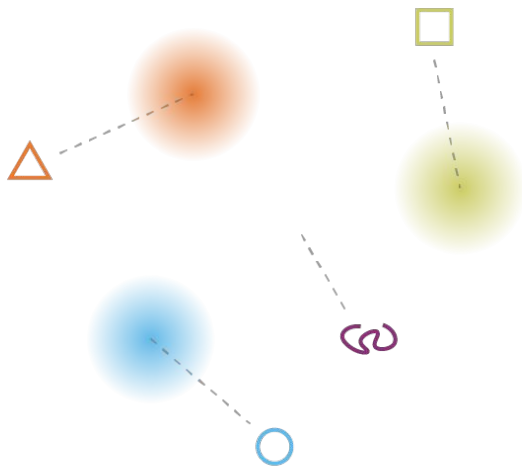


irregular latent space

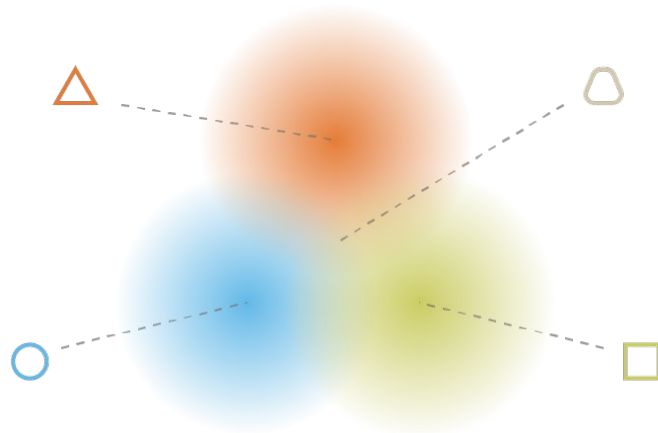


regular latent space





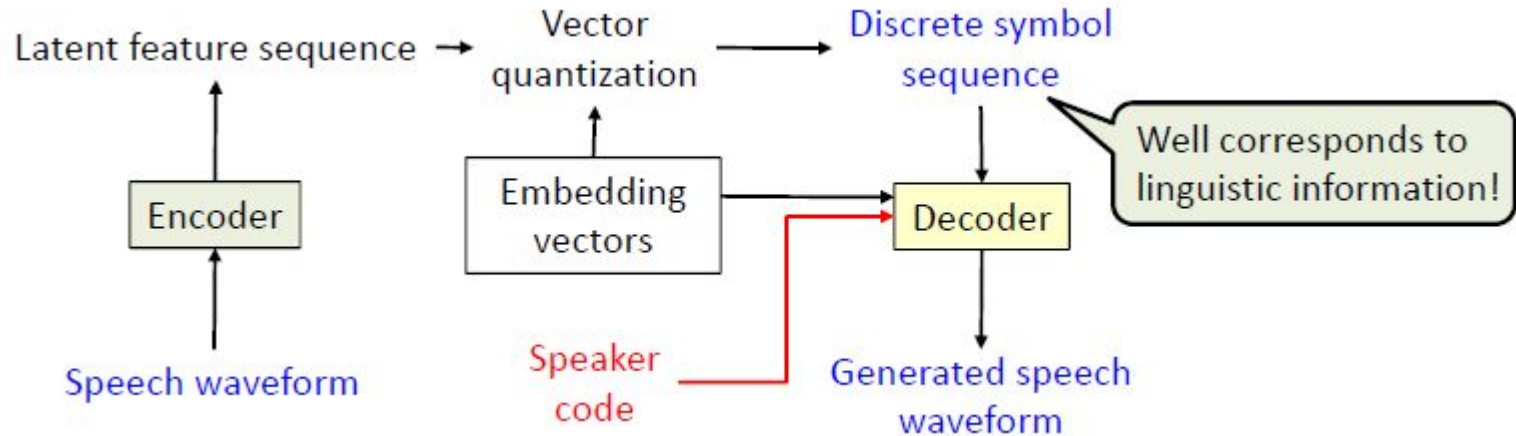
what can happen without regularisation



what we want to obtain with regularisation

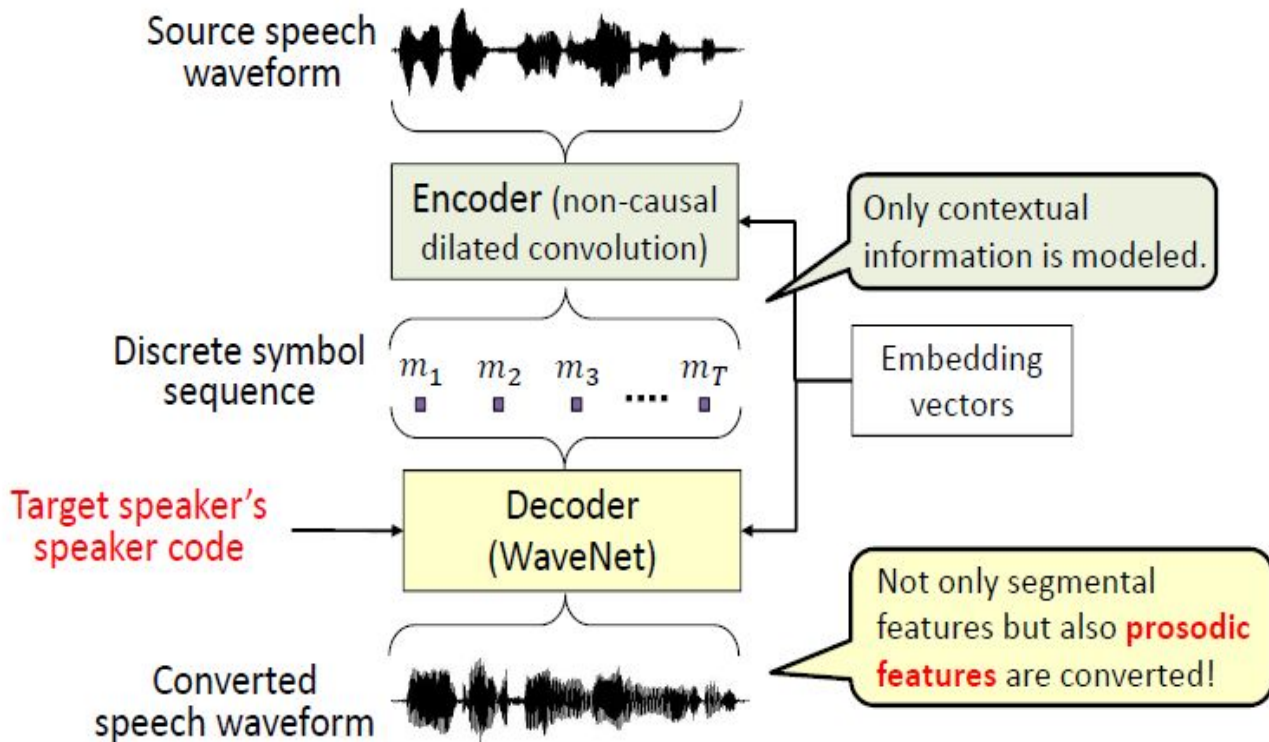
# Vector Quantization VAE (VQ-VAE)

- Directly encode speech waveform into a discrete symbol sequence capturing long-term dependencies (including prosodic features!) by using a dilated convolution network



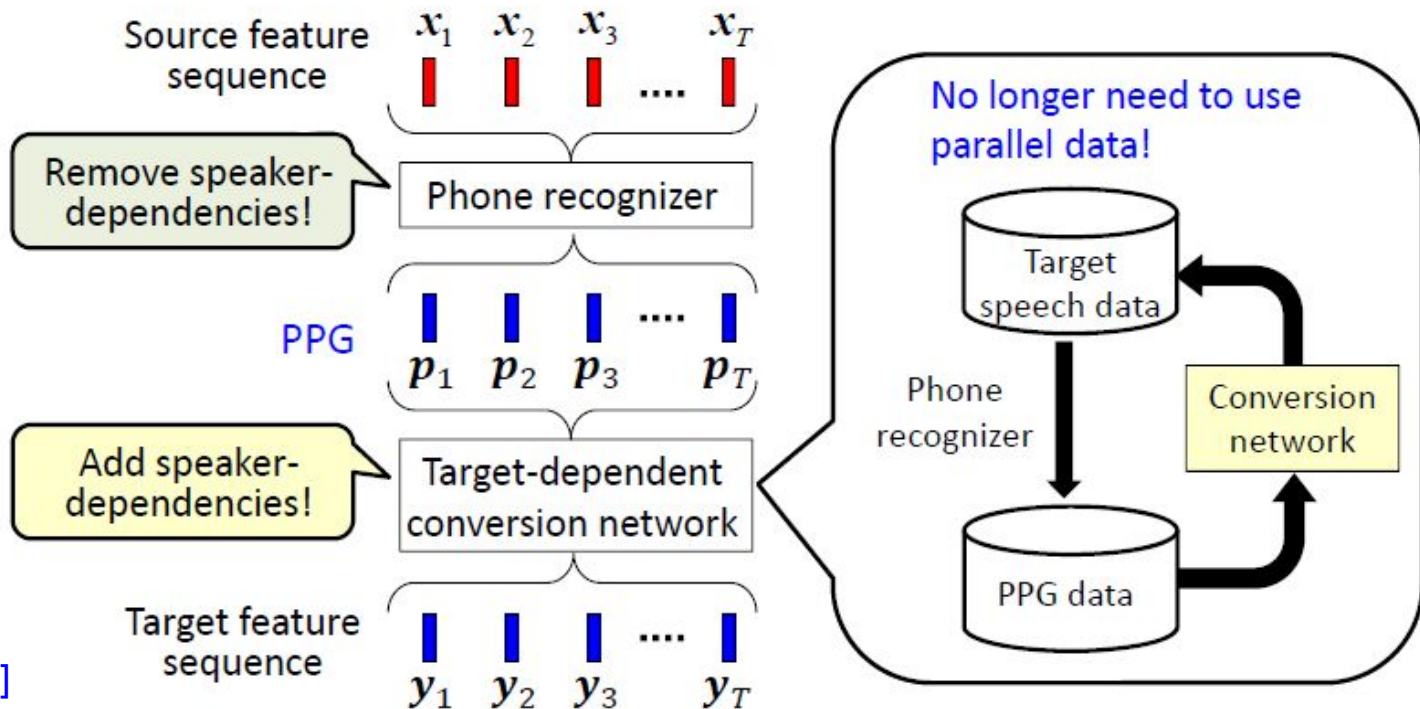
# VC based on VQ-VAE

- Extract phoneme posteriorgram (PPG) as speaker-independent contextual features.



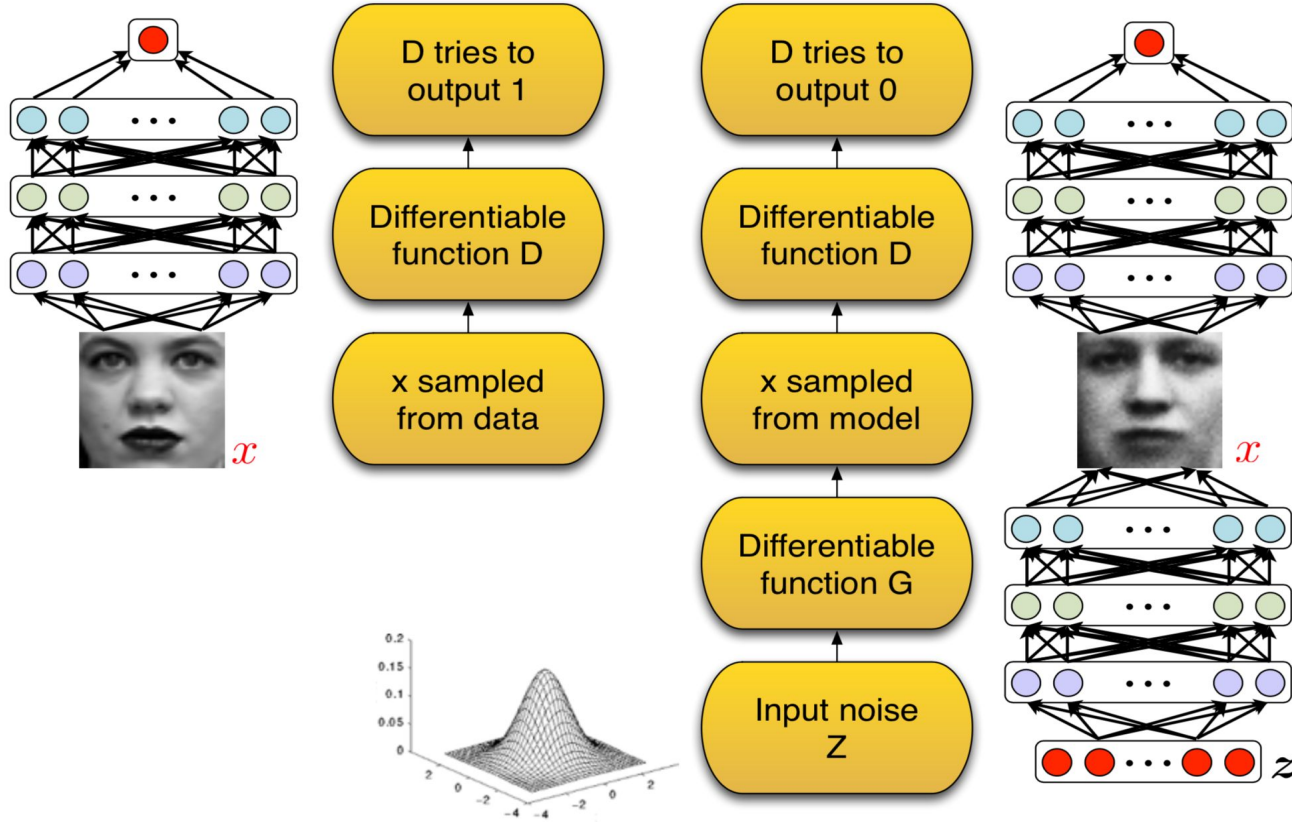
# Phoneme Posteriorgram VC

- Extract phoneme posteriorgram (PPG) as speaker-independent contextual features and use them as input of the conversion network.

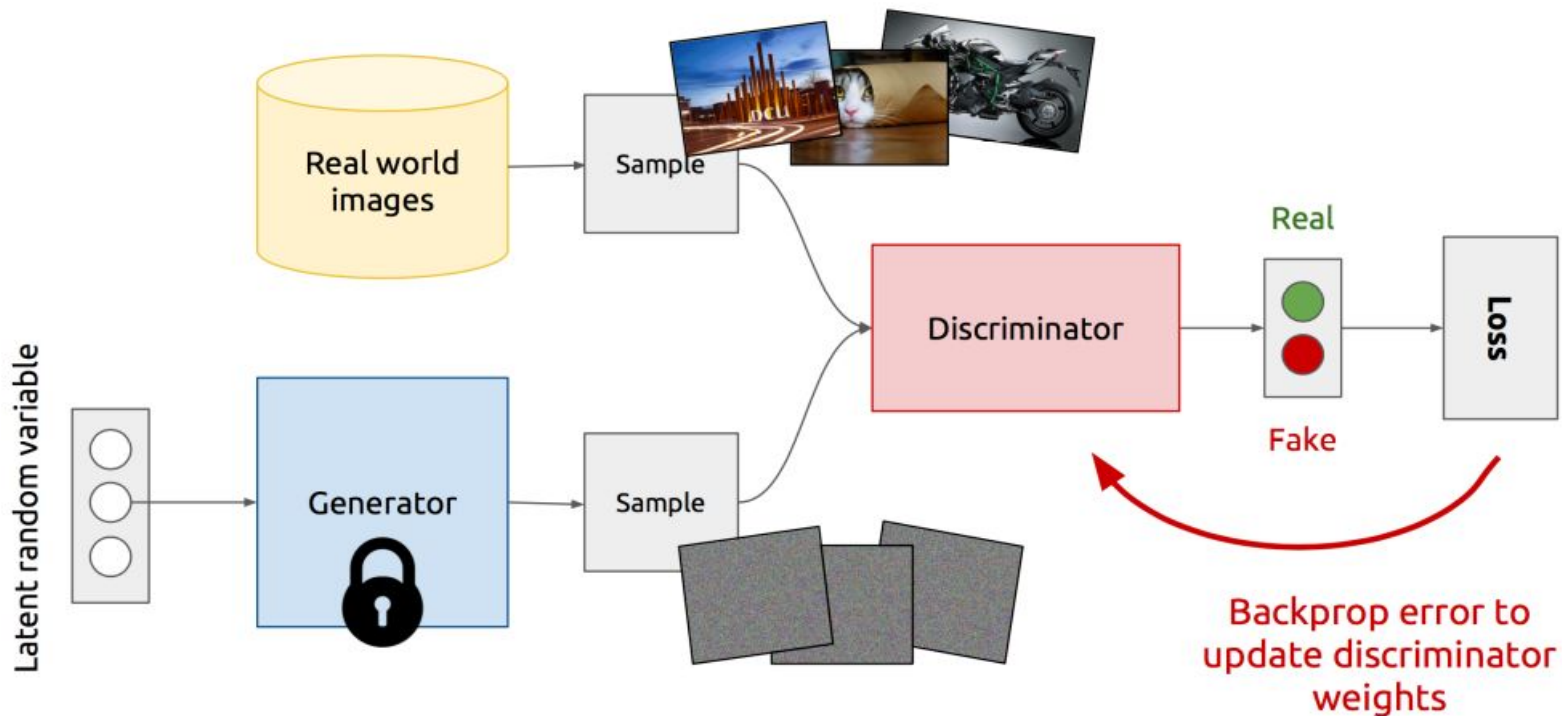


# VC based on Generative Adversarial Networks

# GAN Formulation

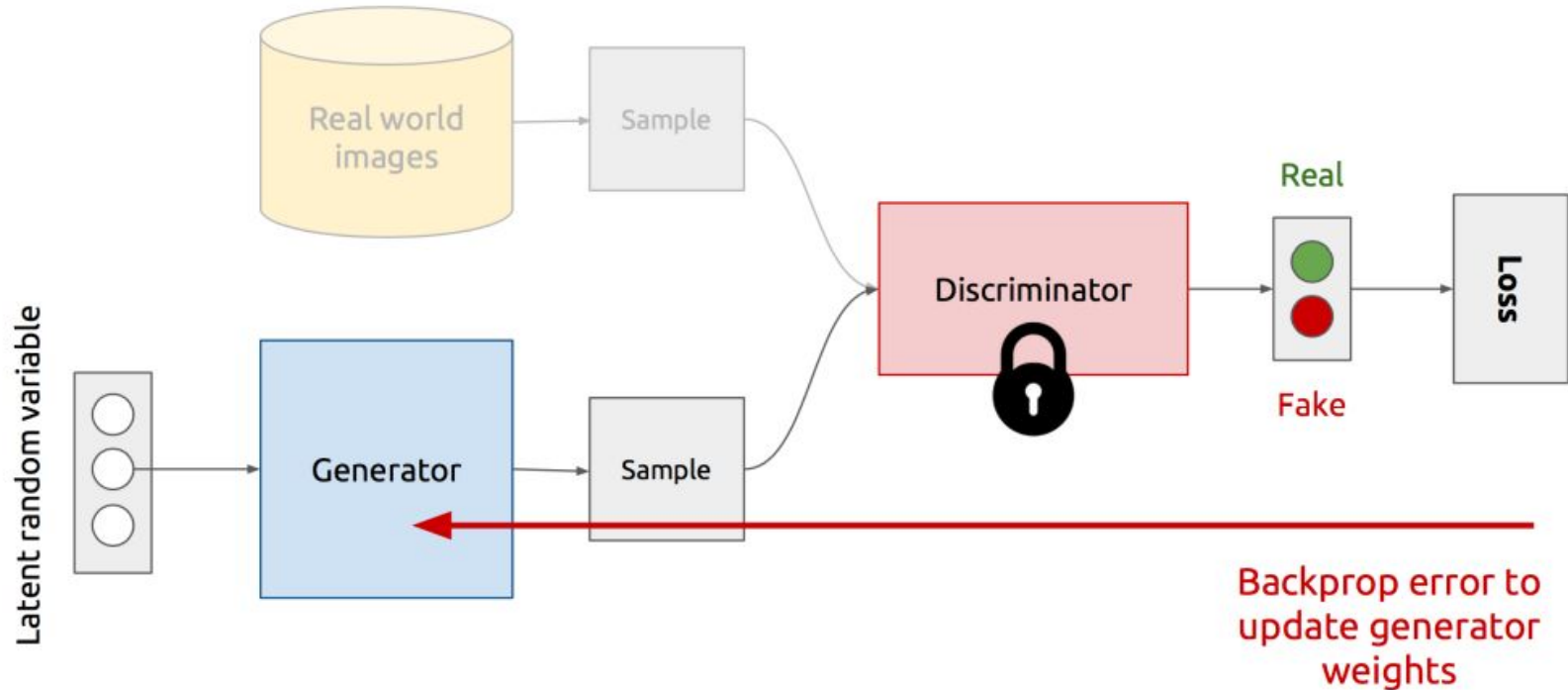


# Discriminator Training





# Generator Training



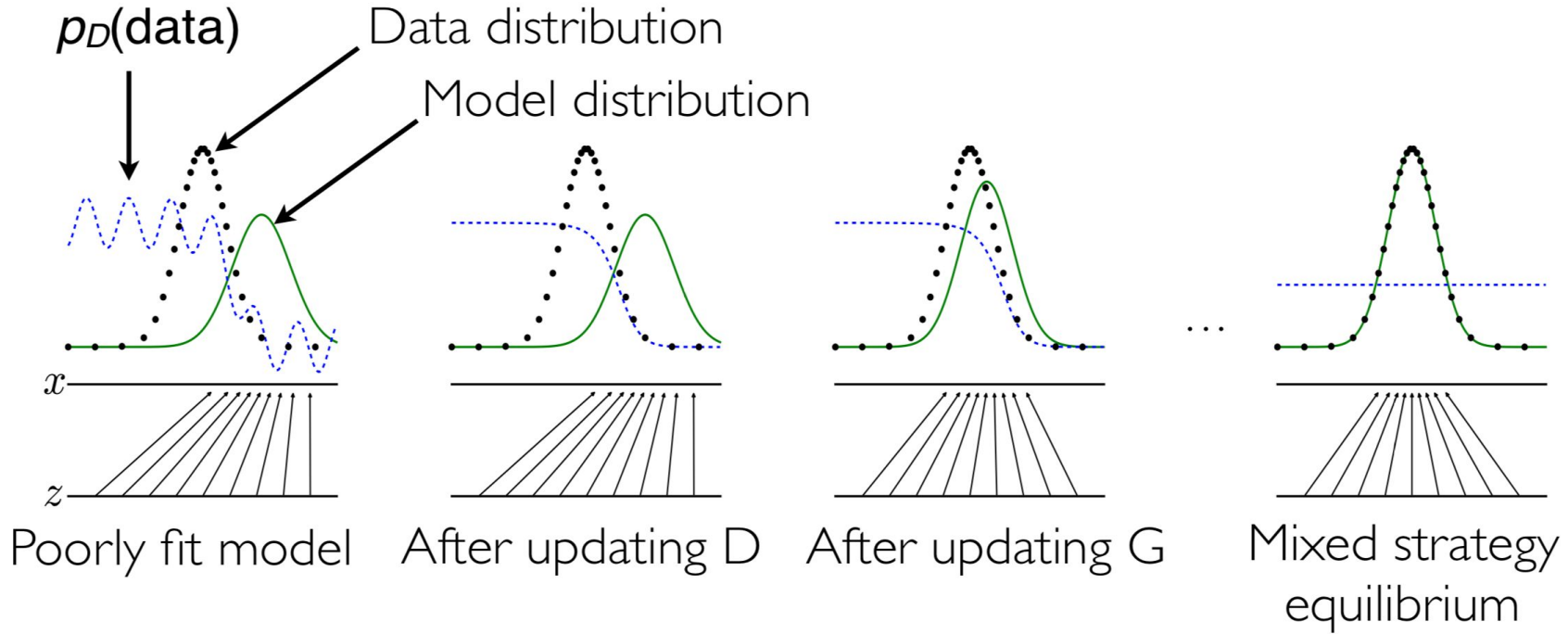
# Mathematical Notations

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

Annotations for the equation above:

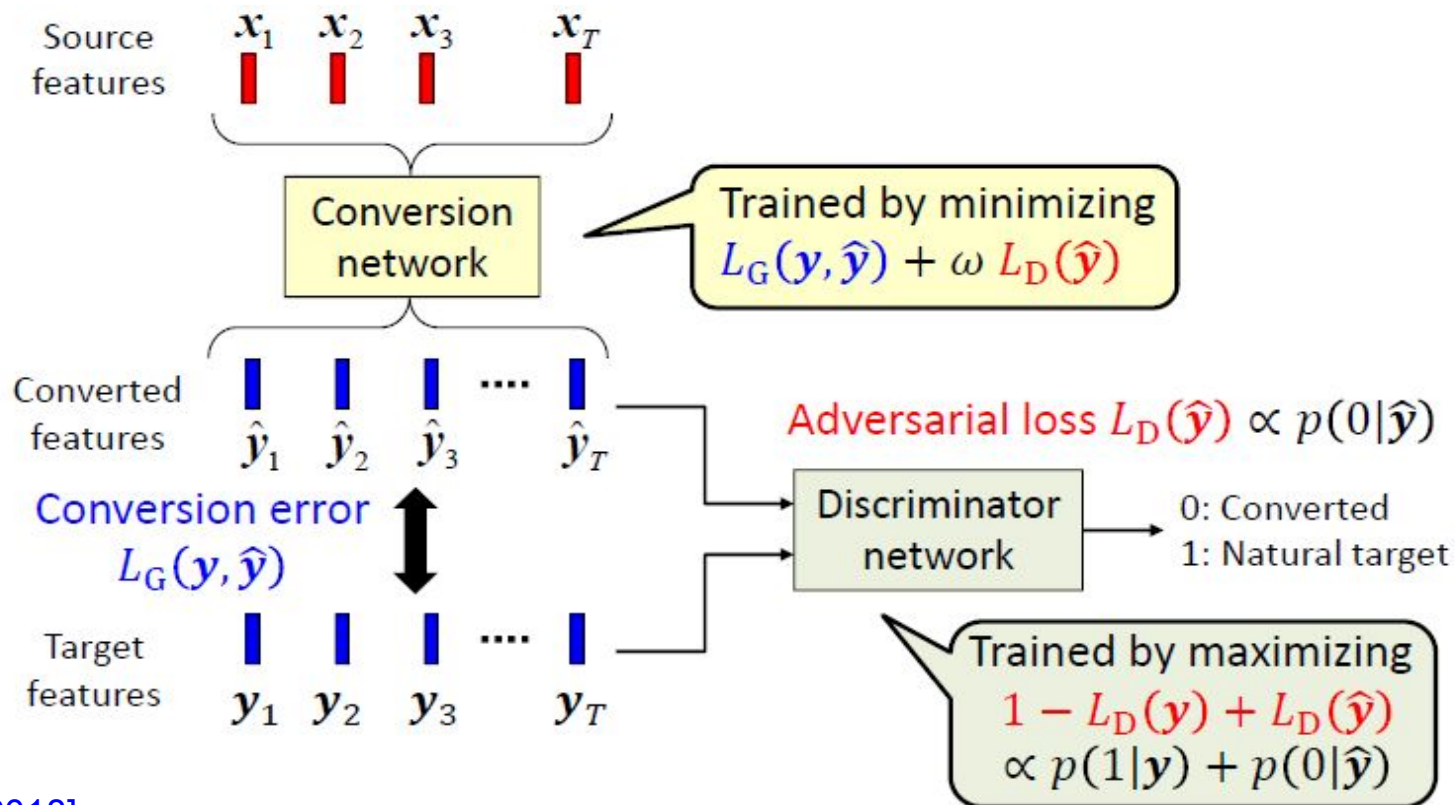
- Value of  $V(D, G)$
- Expectation
- prob. of  $D(\text{real})$
- prob. of  $D(\text{fake})$
- Minimize  $G$
- Maximize  $D$
- $\mathbf{x}$  is sampled from real data
- $\mathbf{z}$  is sampled from  $N(0, I)$
- fake

# Learning GANs



[Goodfellow et al., 2017]

# GAN-based VC

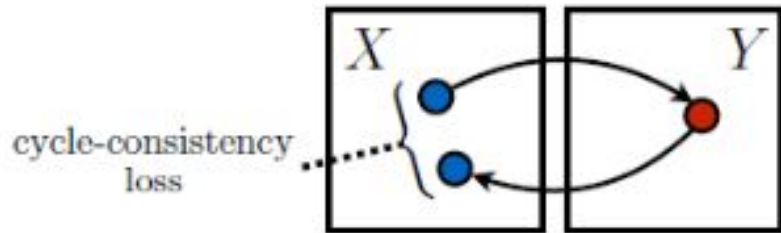
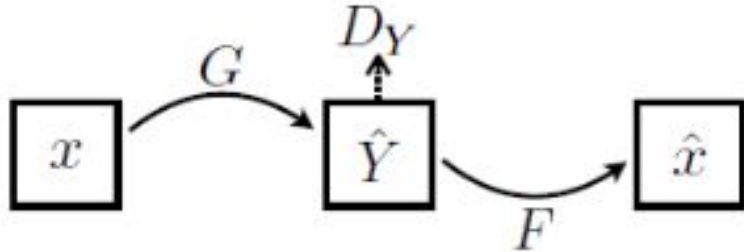


# CycleGAN Voice Conversion

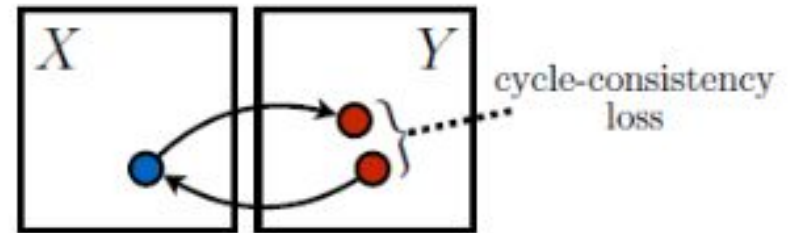
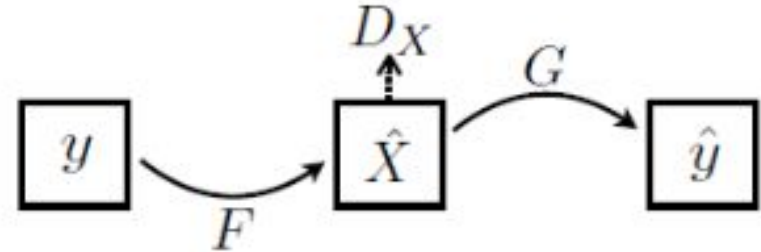
- A non-parallel voice-conversion (VC) method that can learn a mapping from source to target speech without relying on parallel data.
- In a CycleGAN, forward and inverse mappings are simultaneously learned using an adversarial loss and cycle-consistency loss.
- Two important losses are introduced:
  - Adversarial loss
  - cycle-consistency loss
  - identity-mapping loss

# CycleGAN losses

Adversarial loss



Adversarial loss



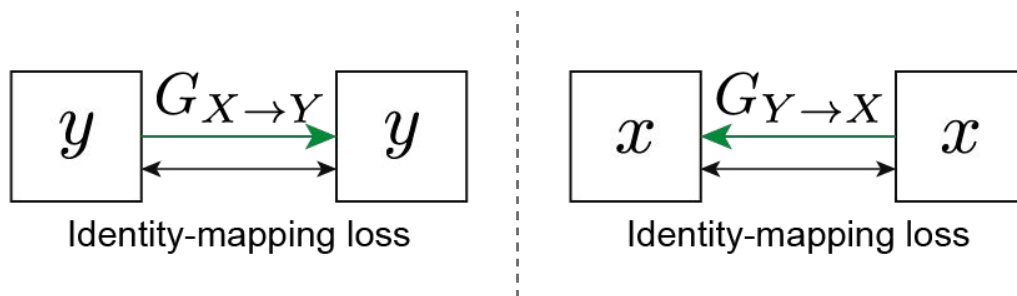
# CycleGAN losses

- Two mapping function (Adversarial loss):  $G$  and  $F$ .  $G : X \rightarrow Y$  and  $F : Y \rightarrow X$
- Cycle-consistency loss:
  - Forward:  $x \rightarrow G(x) \rightarrow F(G(x)) \approx x$
  - Backward:  $y \rightarrow F(y) \rightarrow G(F(y)) \approx y$
- Adversarial loss + cycle-consistency loss:

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

# Identity-mapping loss

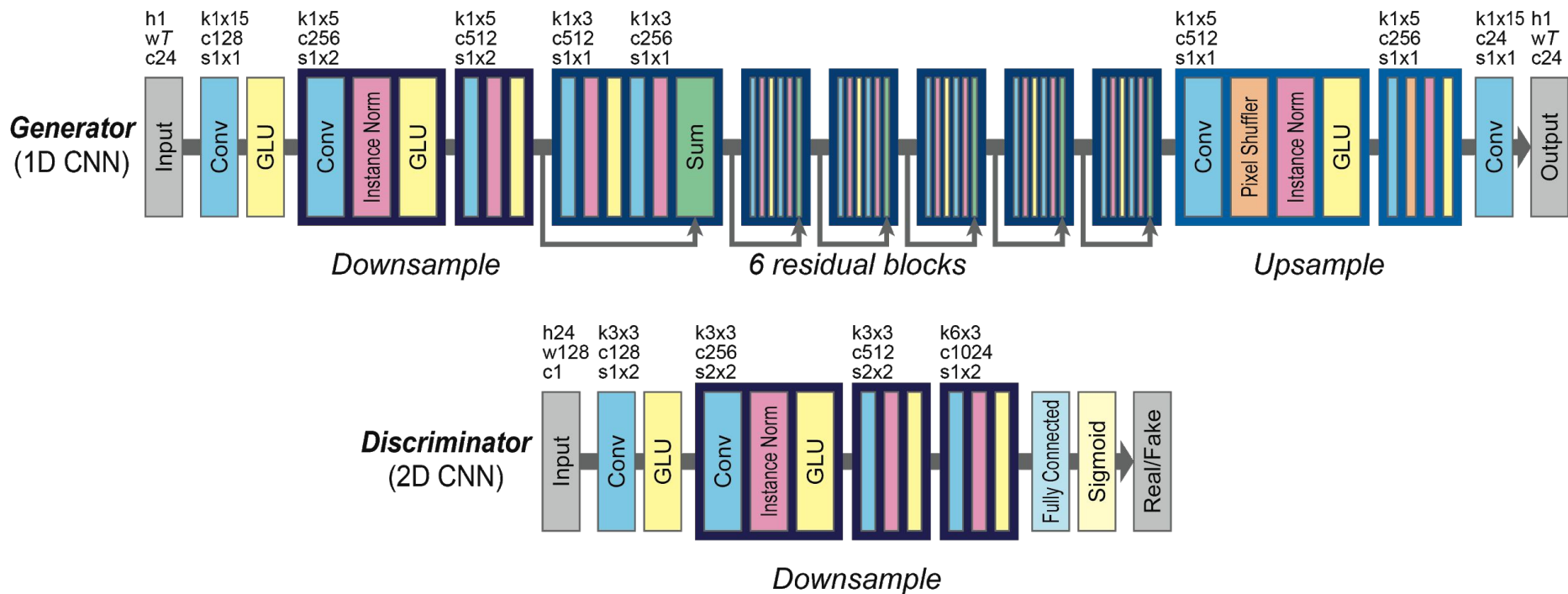
- To encourage linguistic-information preservation, an identity-mapping loss is implemented.
- It encourages the generator to find the mapping that preserves composition between the input and output.



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



# CycleGAN Architecture



# Sound Samples

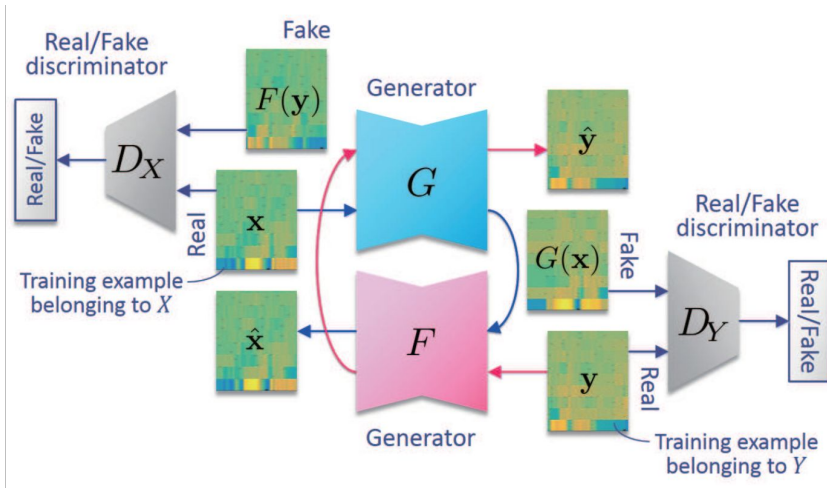
<http://www.kecl.ntt.co.jp/people/kaneko.takuhiro/projects/cyclegan-vc/>

# StarGAN Voice Conversion

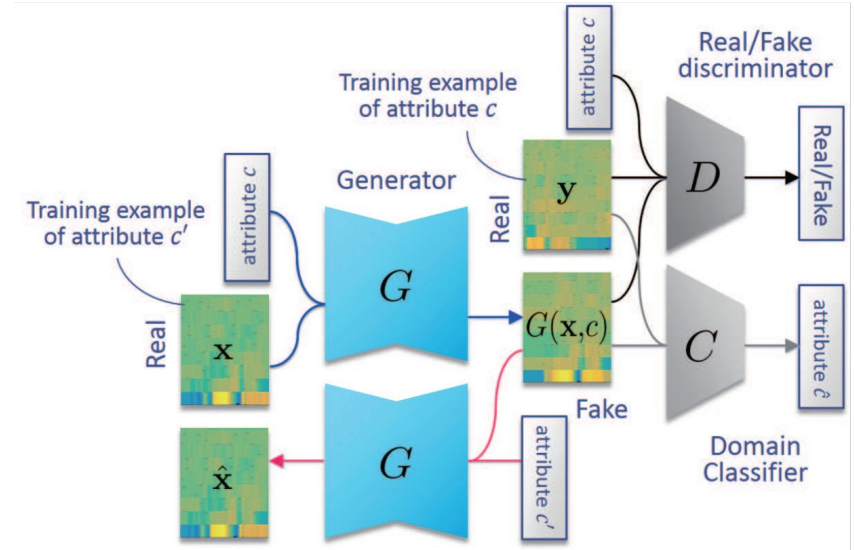
- A non-parallel many-to-many voice conversion (VC) by using a variant of a generative adversarial network called StarGAN.
- Generator (G) takes an acoustic feature with an attribute  $c$  as the inputs and generates an acoustic feature sequence  $y' = G(x, c)$ .
- Discriminator (D) is designed to produce a probability  $D(y, c)$  that an input  $y$  is a real speech feature.
- A domain classifier (C) predicts classes of the input.

# StarGAN training

## CycleGAN



## StarGAN



[Kameoka et. al. 2018]

# StarGAN training losses

## Adversarial loss:

- Adversarial losses for discriminator  $D$  and generator  $G$ , respectively, where  $y$  denotes a training example of an acoustic feature sequence of real speech with attribute  $c$  and  $x$  denotes that with an arbitrary attribute.

$$\begin{aligned}\mathcal{L}_{\text{adv}}^D(D) &= -\mathbb{E}_{c \sim p(c), \mathbf{y} \sim p(\mathbf{y}|c)} [\log D(\mathbf{y}, c)] \\ &\quad - \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}), c \sim p(c)} [\log(1 - D(G(\mathbf{x}, c), c))],\end{aligned}$$

$$\mathcal{L}_{\text{adv}}^G(G) = -\mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}), c \sim p(c)} [\log D(G(\mathbf{x}, c), c)],$$

# StarGAN training losses

## Domain Classification loss:

- Domain classification losses for classifier C and generator G is described.

$$\mathcal{L}_{\text{cls}}^C(C) = - \mathbb{E}_{c \sim p(c), \mathbf{y} \sim p(\mathbf{y}|c)} [\log p_C(c|\mathbf{y})],$$

$$\mathcal{L}_{\text{cls}}^G(G) = - \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}), c \sim p(c)} [\log p_C(c|G(\mathbf{x}, c))],$$

# StarGAN training losses

## Cycle Consistency Loss:

- To encourage  $G(x, c)$  to be a bijection, a cycle consistency loss is implemented, where  $x$  denotes an acoustic feature sequence of real speech with attribute  $c'$ .

$$\mathcal{L}_{\text{cyc}}(G) = \mathbb{E}_{c' \sim p(c), \mathbf{x} \sim p(\mathbf{x}|c'), c \sim p(c)} [\|G(G(\mathbf{x}, c), c') - \mathbf{x}\|_{\rho}],$$

# StarGAN training losses

## Identity mapping loss:

- Ensure that an input into  $G$  will remain unchanged when the input already belongs to the target attribute  $c'$ .

$$\mathcal{L}_{\text{id}}(G) = \mathbb{E}_{c' \sim p(c), \mathbf{x} \sim p(\mathbf{x}|c')} [\|G(\mathbf{x}, c') - \mathbf{x}\|_{\rho}],$$



# StarGAN Objective Function

## Objective function :

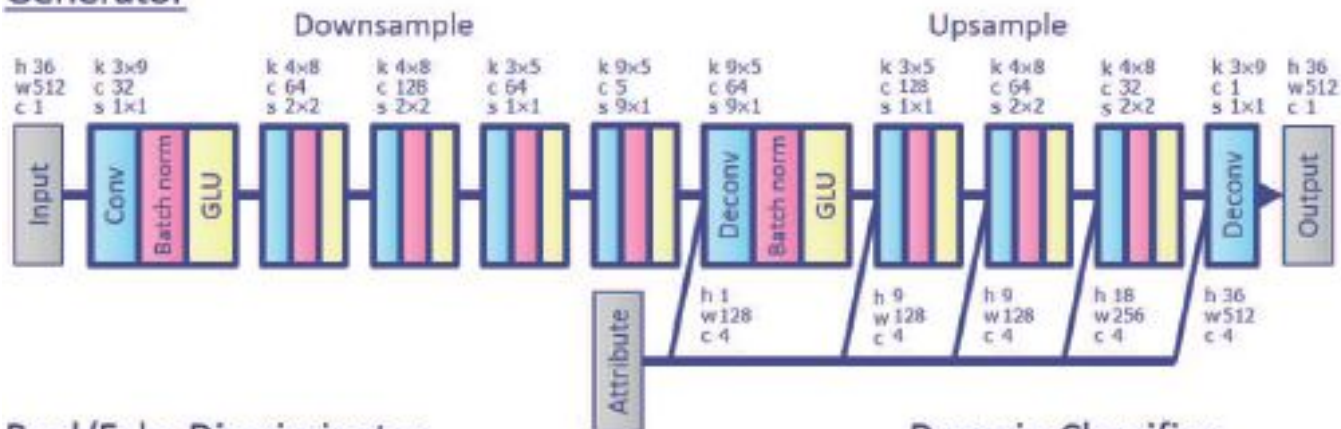
- The full objectives of StarGAN-VC to be minimized with respect to  $G$ ,  $D$  and  $C$  are

$$\mathcal{I}_G(G) = \mathcal{L}_{\text{adv}}^G(G) + \lambda_{\text{cls}} \mathcal{L}_{\text{cls}}^G(G) + \lambda_{\text{cyc}} \mathcal{L}_{\text{cyc}}(G) + \lambda_{\text{id}} \mathcal{L}_{\text{id}}(G)$$

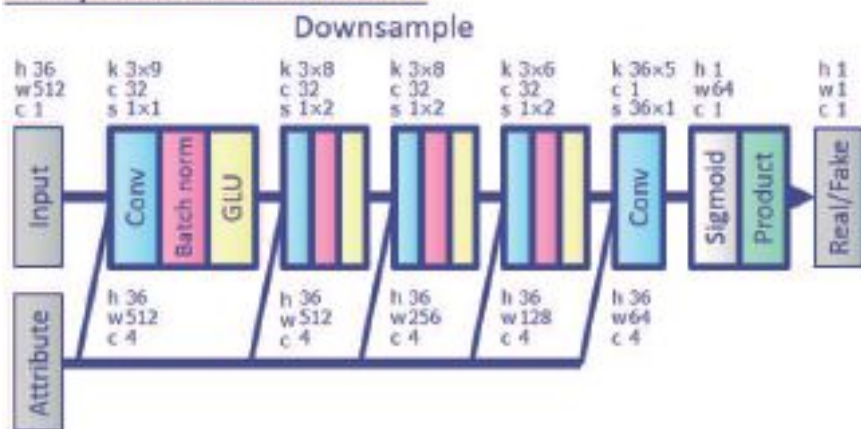
$$\mathcal{I}_D(D) = \mathcal{L}_{\text{adv}}^D(D),$$

$$\mathcal{I}_C(C) = \mathcal{L}_{\text{cls}}^C(C),$$

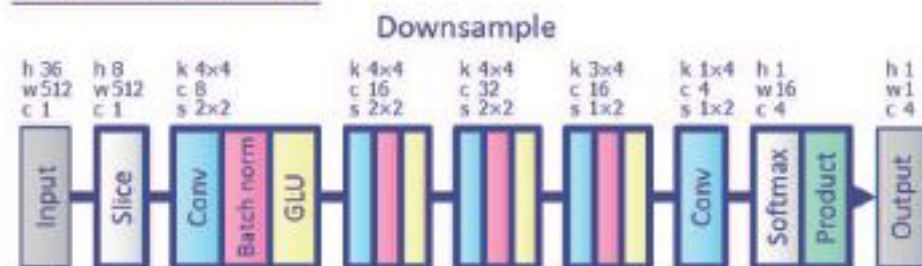
## Generator



## Real/Fake Discriminator

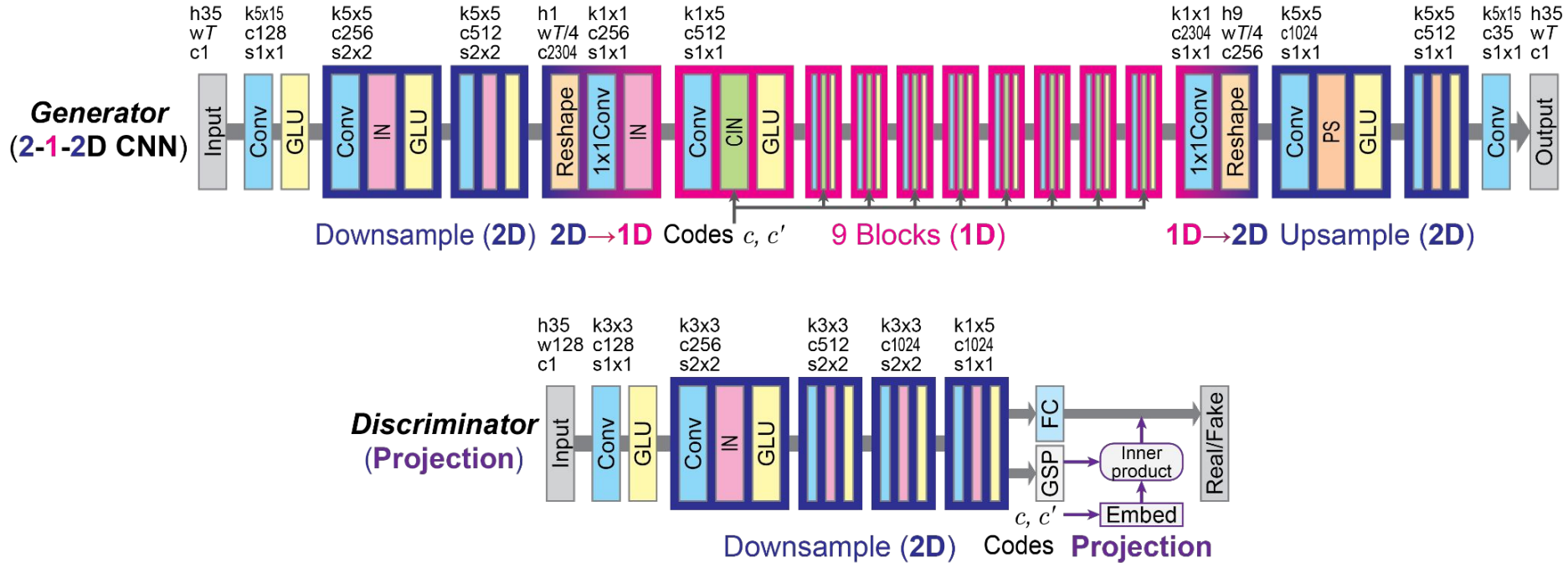


## Domain Classifier



[Kameoka et. al. 2018]

# Modified StarGAN



# Rethinking Conditional Methods

- source-and-target conditional adversarial loss defined as

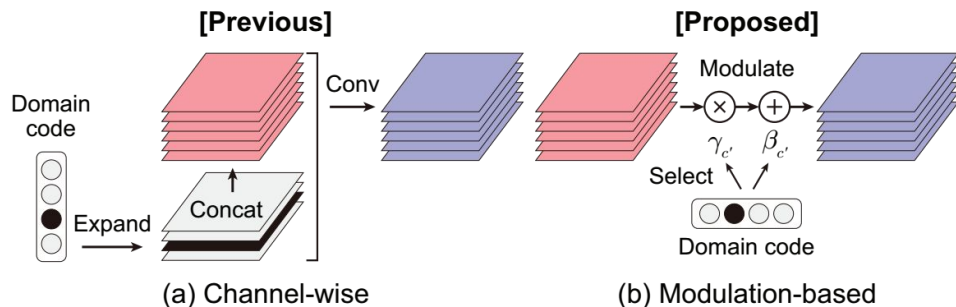
$$\begin{aligned}\mathcal{L}_{st-adv} = & \mathbb{E}_{(\mathbf{x},c) \sim P(\mathbf{x},c), c' \sim P(c')} [\log D(\mathbf{x}, c', c)] \\ & + \mathbb{E}_{(\mathbf{x},c) \sim P(\mathbf{x},c), c' \sim P(c')} [\log D(G(\mathbf{x}, c, c'), c, c')],\end{aligned}$$

# Rethinking Conditional Methods

- Given the feature  $\mathbf{f}$ , conditional instance normalization (CIN) conducts the following procedure:

$$\text{CIN}(\mathbf{f}; c') = \gamma_{c'} \left( \frac{\mathbf{f} - \mu(\mathbf{f})}{\sigma(\mathbf{f})} \right) + \beta_{c'},$$

where  $\mu(\mathbf{f})$  and  $\sigma(\mathbf{f})$  are the average and standard deviation of  $\mathbf{f}$  that are calculated over for each instance.  $\gamma_{c'}$  and  $\beta_{c'}$  are domain-specific scale and bias parameters that allow the modulation to be transformed in a domain-specific manner.



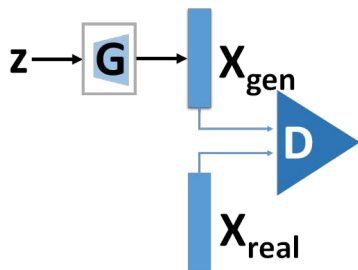
# Sound Samples

<http://www.kecl.ntt.co.jp/people/kameoka.hirokazu/Demos/stargan-vc/>

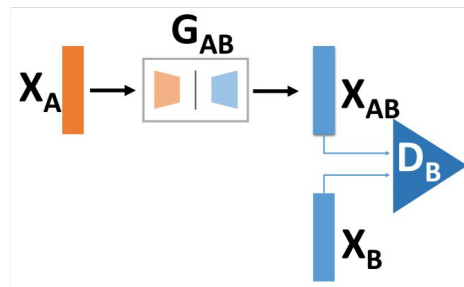
<http://www.kecl.ntt.co.jp/people/kaneko.takuhiro/projects/stargan-vc2/index.html>

# Style Conversion (VoiceGAN)

- Voice style impersonation, where one person attempts to mimic the voice of another to sound like the other person, is a complex phenomenon.



The original GAN model



Style transfer by GAN

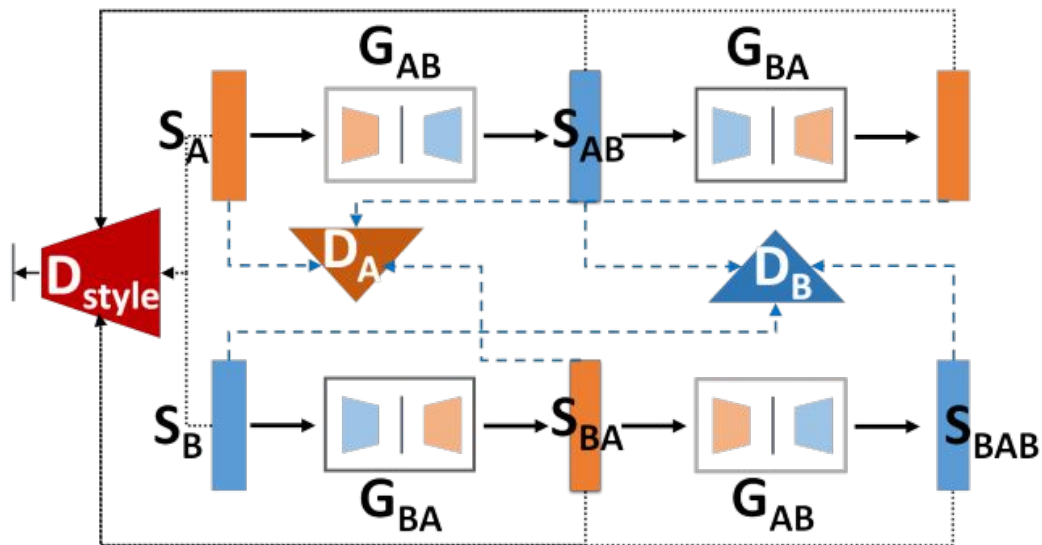
Mathematical notations:

$$L_G = E_{x_A \sim P_A} [\log(1 - D_B(x_{AB}))]$$

$$L_D = -E_{x_B \sim P_B} [\log D_B(x_B)] - E_{x_A \sim P_A} [\log(1 - D_B(x_{AB}))]$$

[Gao et al., 2018]

# Style Conversion (VoiceGAN)



- $D_A$  and  $D_B$  discriminate between real and fake data.
- $G_{AB}$  transforms for style A to style B, whereas  $G_{BA}$  is the opposite.
- The discriminator  $D_{style}$  determines if the original and transformed signals match the desired style.



# Style Conversion (VoiceGAN)

- Training objectives to be minimized for the generator and discriminator are represented by  $L_G$  and  $L_D$  respectively as follows:

$$L_G = L_{GAN_{AB}} + L_{GAN_{BA}} = L_{G_B} + L_{CONST_A} + L_{G_A} + L_{CONST_B}$$

$$L_D = L_{D_A} + L_{D_B} + L_{D_{STYLE}}$$

- Reconstruction loss:

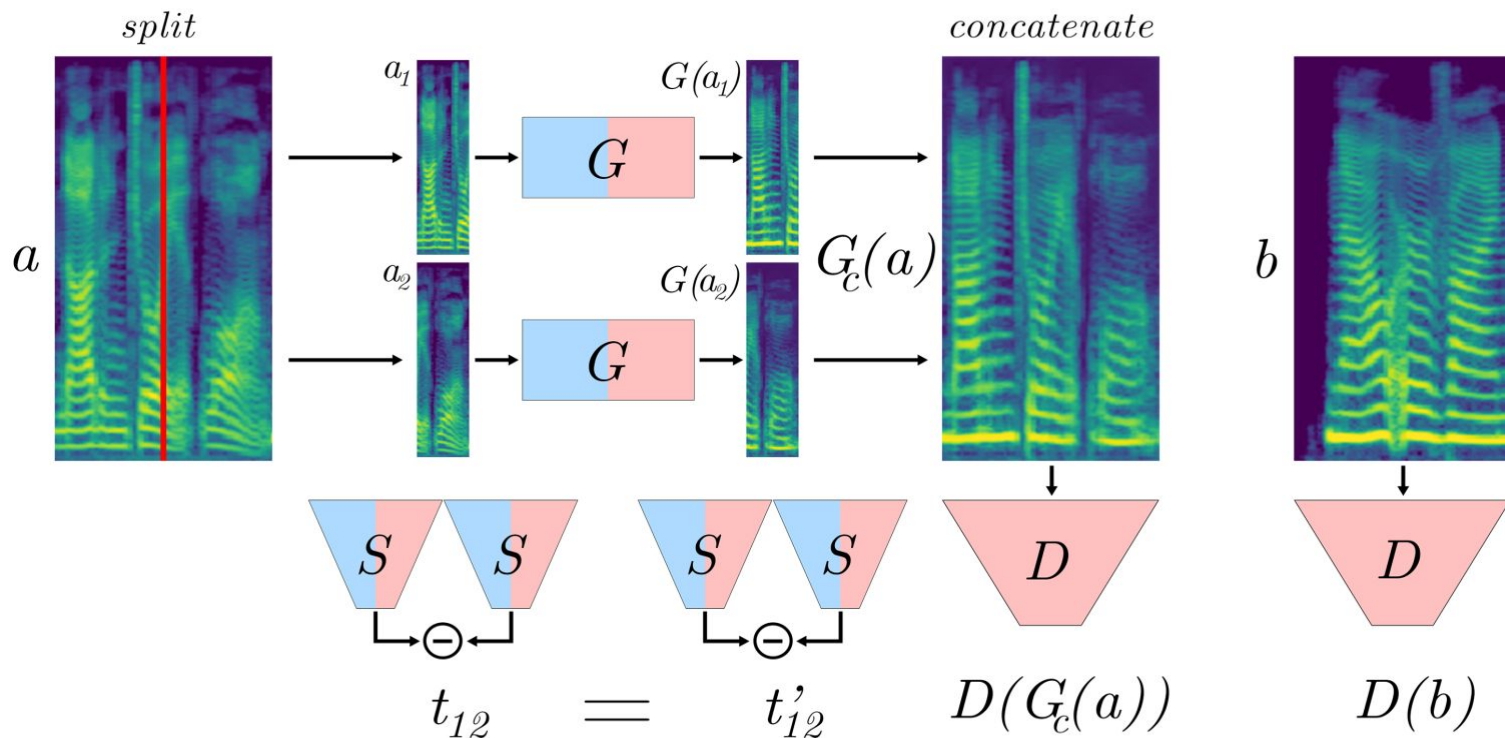
$$L_{CONST_A} = d(G_{BA}(G_{AB}(x_A)), x_A)$$

- The discriminator  $D_S$  determines if the original and transformed signals match the desired style:

$$L_{D_{STYLE}} = L_{D_{STYLE-A}} + L_{D_{STYLE-B}}$$

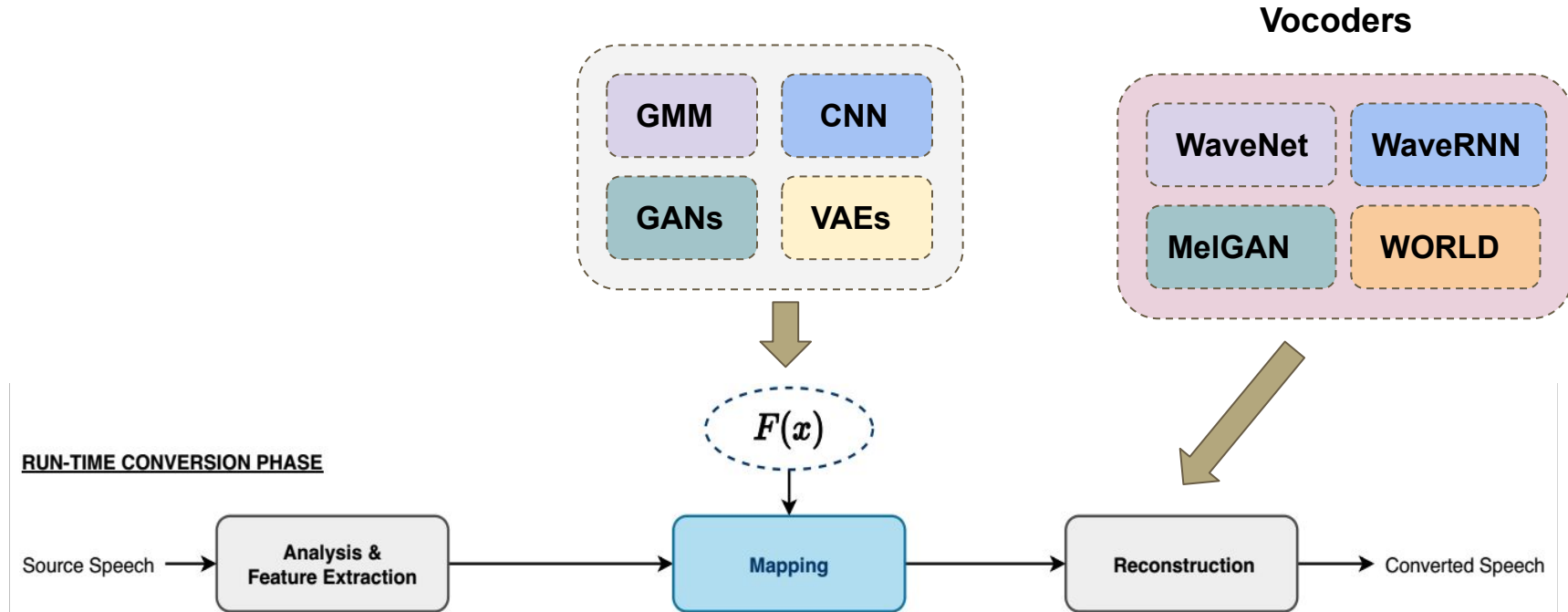
$$L_{D_{STYLE-A}} = d(D_S(x_A), label_A) + d(D_S(x_{AB}), label_B) + d(D_S(x_{ABA}), label_A)$$

# MeLGAN VC

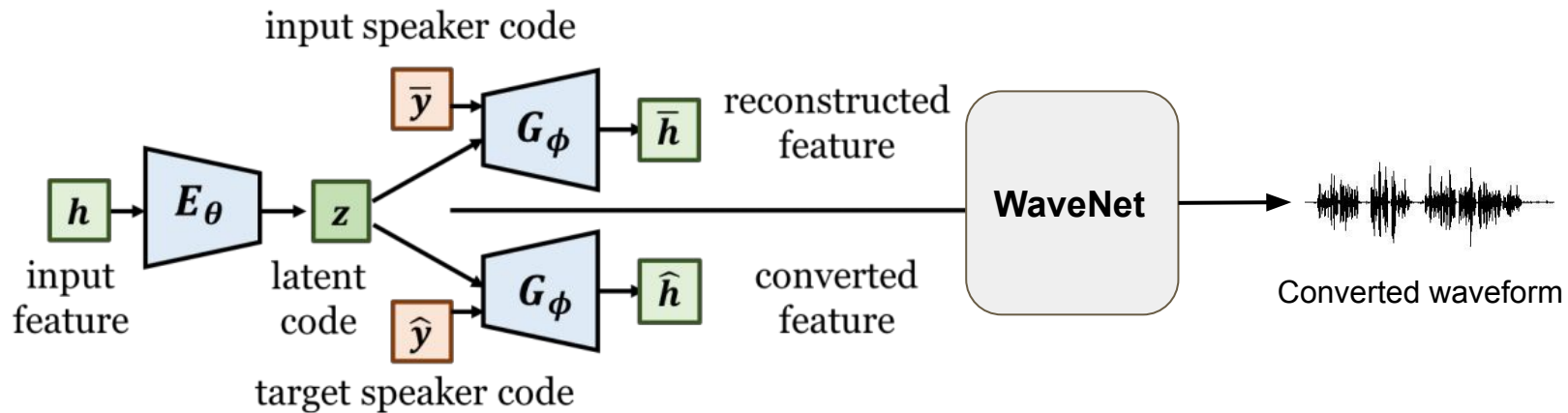


# Various Vocoders in VC

# General Framework



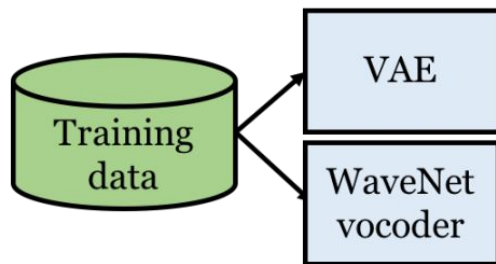
# WaveNet Vocoder in VAE-VC



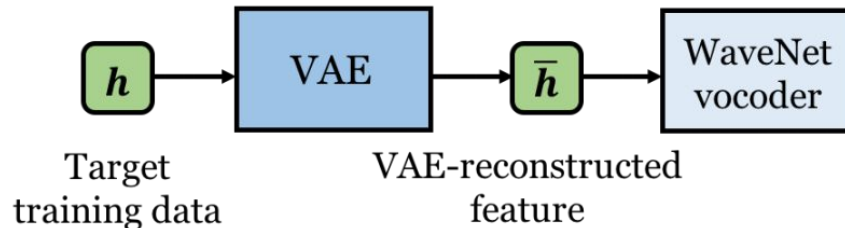
A general framework of WaveNet vocoder in voice conversion.

[Huang et. al. 2019]

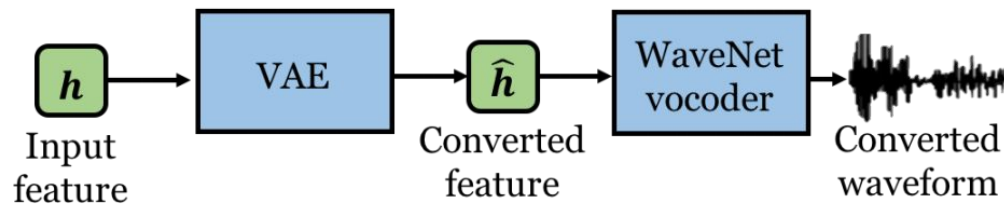
# Training Protocol



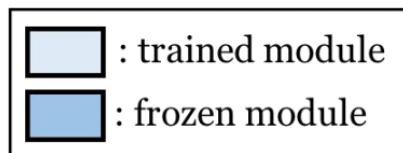
**Step 1: VAE and WaveNet vocoder training phase**



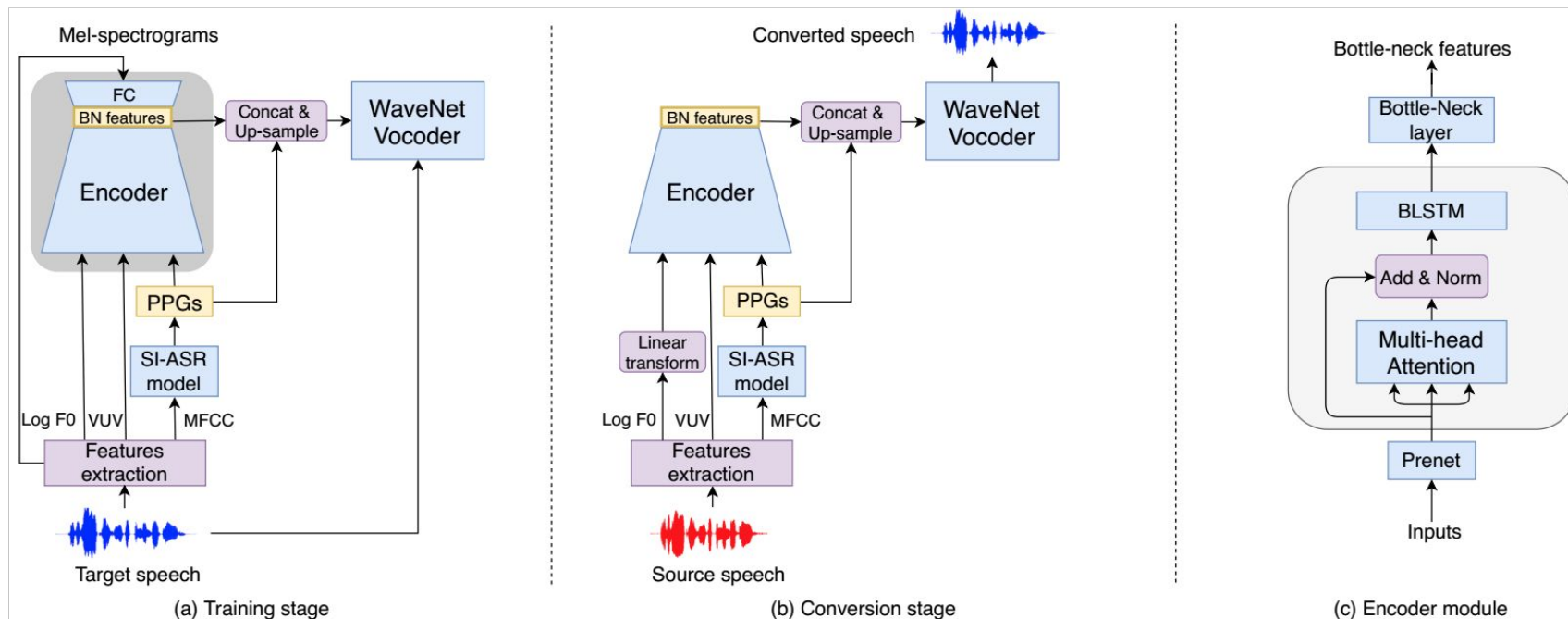
**Step 2: WaveNet vocoder training phase**



**Step 3: Conversion phase**

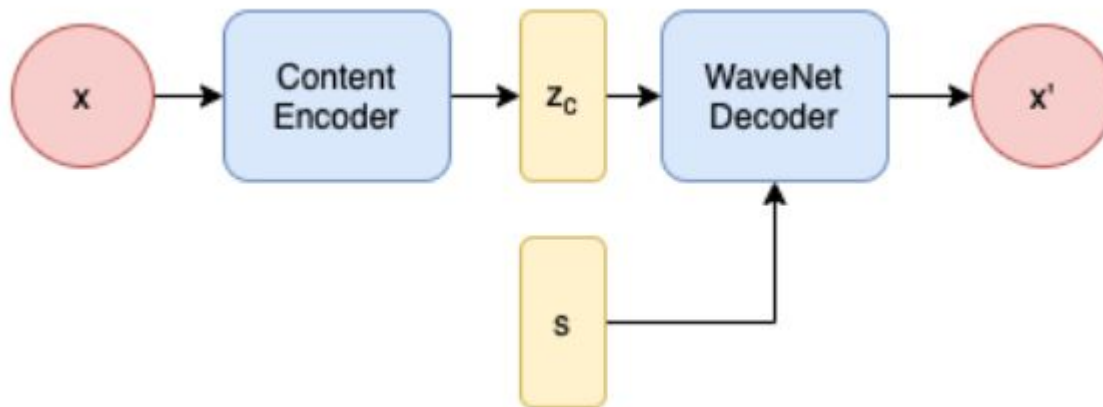


# Jointly Trained Conversion Model and Vocoder



# WaveNet Auto-encoders

- WaveNet is used as the decoder and to generate waveform data directly from the latent representation.



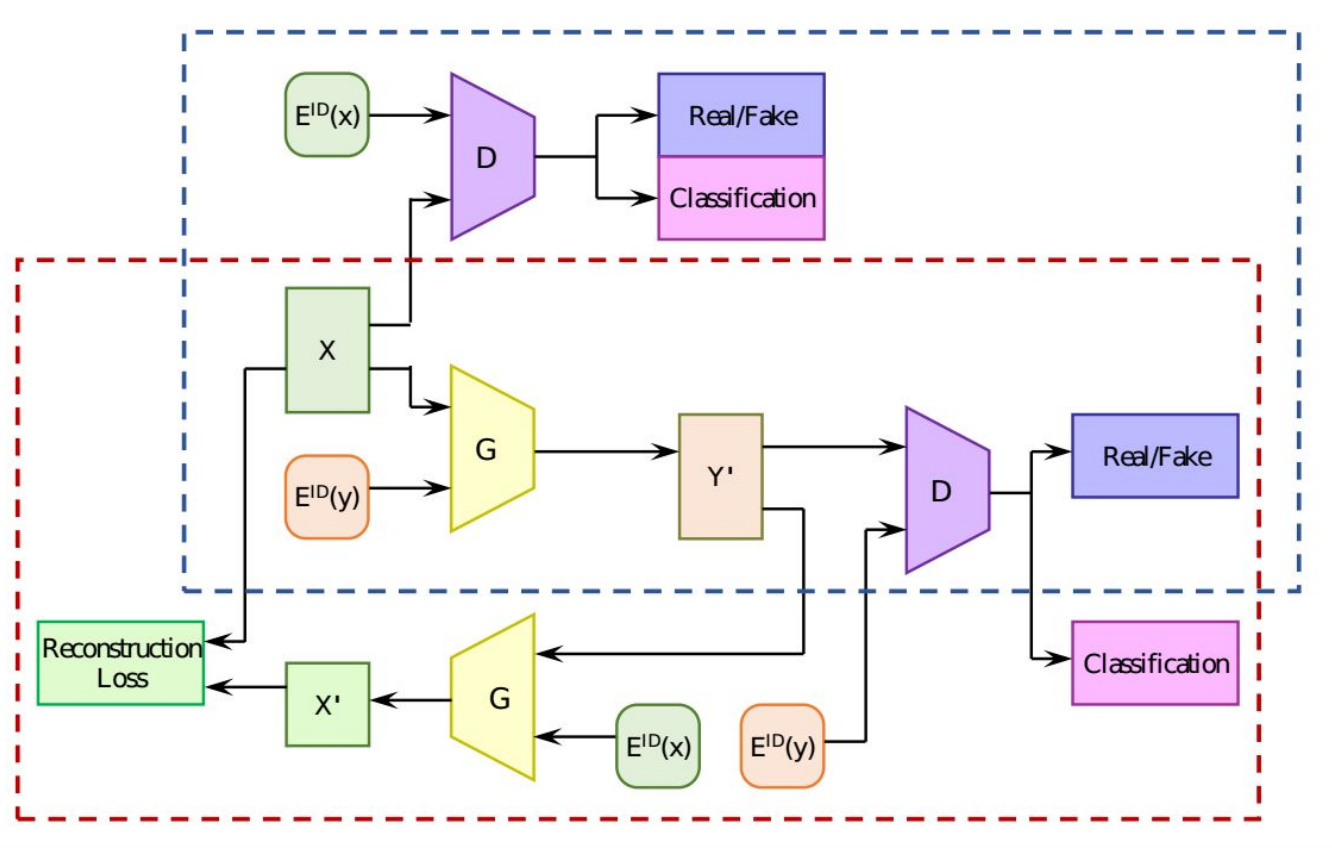


# One-shot Voice Conversion

# One-shot VC

- The target speaker is unseen in training dataset or both source and target speakers are unseen in the training dataset.
- 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 would be a weighted combination of the timbres the speakers seen in the dataset.

# One-shot StarGAN VC

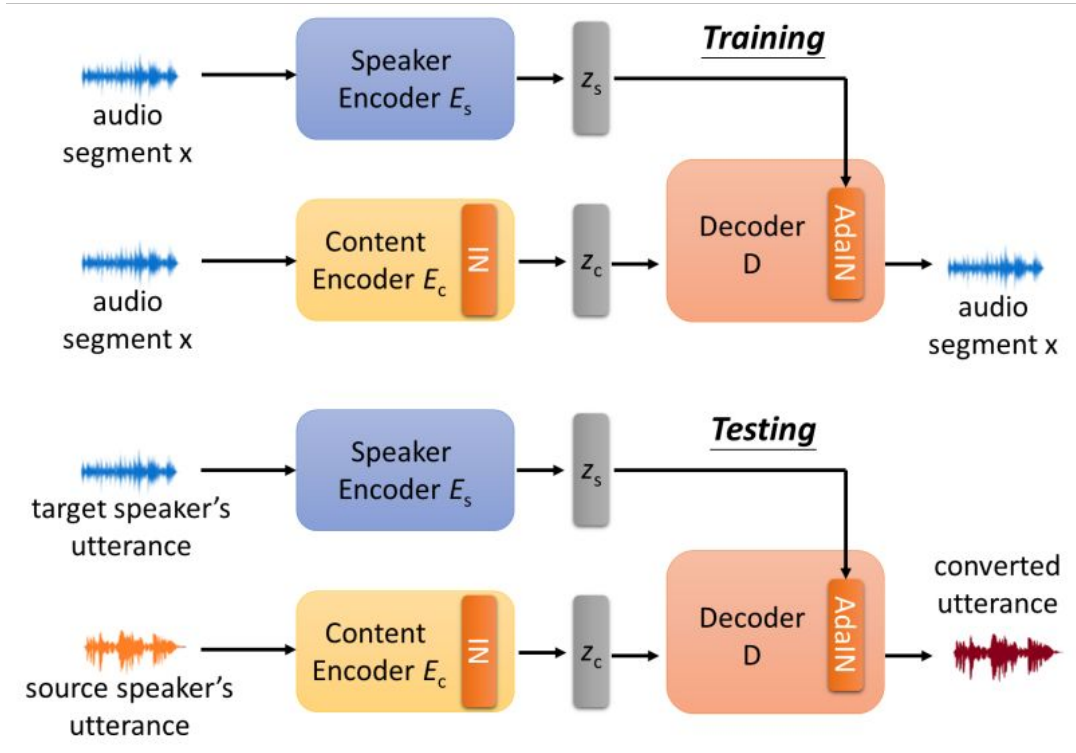


[Wang et. al. 2020]

# Representations Learning

- An utterance can be factorized into a speaker plus a content representation.
- To disentangle speaker and content representation, three components is employed : a speaker encoder, a content encoder and a decoder
- The speaker encoder is trained to encode the speaker information.
- The content encoder is trained to encode only the linguistic information.
- The task of the decoder is to synthesize the voice back by combining these two representations.

# Representations Learning



- $E_s$  is speaker encoder
- $E_c$  is content encoder
- $D$  is decoder.
- $IN$  is instance normalization
- $AdaIN$  represents adaptive instance normalization layer.

[Chou et. al. 2019]

# Representations Learning

- The objective function for VAE training

$$\min_{\theta_{E_S}, \theta_{E_C}, \theta_D} L(\theta_{E_S}, \theta_{E_C}, \theta_D) = \lambda_{rec} L_{rec} + \lambda_{kl} L_{kl}$$

- The reconstruction loss is given as

$$L_{rec}(\theta_{E_S}, \theta_{E_C}, \theta_D) = \mathbb{E}_{x \sim p(x), z_c \sim p(z_c|x)} [\|D(E_S(x), z_c) - x\|_1^1].$$

- The divergence term is given as in

$$L_{kl}(\theta_{E_C}) = \mathbb{E}_{x \sim p(x)} [\|E_C(x)\|_2^2].$$

# Bibliography

- I. Goodfellow et al. "Generative adversarial nets," in Proc. NIPS, 2014.
- Arjovsky et al. "Wasserstein GAN," arXiv preprint arXiv:1701.07875, 2017.
- I. Gulrajani et al., "Improved training of Wasserstein GANs," in Advances in Neural Information Processing Systems, 2017, pp. 5767–5777.
- Hung-yi Lee, and Yu Tsao. "Generative Adversarial Network and its Applications to Speech Signal and Natural Language Processing." tutorial in ICASSP, 2018.
- Hsu, Chin-Cheng, Hsin-Te Hwang, Yi-Chiao Wu, Yu Tsao, and Hsin-Min Wang. "Voice conversion from non-parallel corpora using variational auto-encoder." In *2016 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA)*, pp. 1-6. IEEE, 2016.
- A. Kain and M. W. Macon. Spectral voice conversion for text-to-speech synthesis. Proc. IEEE ICASSP, pp. 285–288, 1998.
- T. Kaneko, H. Kameoka. CycleGAN-VC: non-parallel voice conversion using cycle-consistent adversarial networks. Proc. EUSIPCO, pp. 2114–2118, 2018.
- H. Kameoka, T. Kaneko, K. Tanaka, N. Hojo. StarGAN-VC: non-parallel many-to-many voice conversion using star generative adversarial networks. Proc. IEEE SLT, pp. 266–272, 2018.
- Kaneko, T., Kameoka, H., Tanaka, K. and Hojo, N., 2019. StarGAN-VC2: Rethinking Conditional Methods for StarGAN-Based Voice Conversion}}. Proc. Interspeech 2019, pp. 679-683.

# Bibliography

- Gao, Yang, Rita Singh, and Bhiksha Raj. "Voice impersonation using generative adversarial networks." In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 2506-2510. IEEE, 2018.
- Huang, Wen-Chin, Yi-Chiao Wu, Hsin-Te Hwang, Patrick Lumban Tobing, Tomoki Hayashi, Kazuhiro Kobayashi, Tomoki Toda, Yu Tsao, and Hsin-Min Wang. "Refined wavenet vocoder for variational autoencoder based voice conversion." In *2019 27th European Signal Processing Conference (EUSIPCO)*, pp. 1-5. IEEE, 2019.
- Liu, Songxiang, Yüewen Cao, Xixin Wu, Lifa Sun, Xunying Liu, and Helen Meng. "Jointly Trained Conversion Model and WaveNet Vocoder for Non-Parallel Voice Conversion Using Mel-Spectrograms and Phonetic Posteriorgrams." In *INTERSPEECH*, pp. 714-718. 2019.
- Pasini, Marco. "Melgan-vc: Voice conversion and audio style transfer on arbitrarily long samples using spectrograms." *arXiv preprint arXiv:1910.03713* (2020).
- Chen, Mingjie, and Thomas Hain. "Unsupervised acoustic unit representation learning for voice conversion using wavenet auto-encoders." *arXiv preprint arXiv:2008.06892* (2020).
- 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.
- Chou, Ju-chieh, and Hung-Yi Lee. "One-Shot Voice Conversion by Separating Speaker and Content Representations with Instance Normalization}." *Proc. Interspeech 2019* (2019): 664-668.



# Bibliography

- M. Abe, S. Nakamura, K. Shikano, H. Kuwabara. Voice conversion through vector quantization. J. Acoust. Soc. Jpn (E), Vol. 11, No. 2, pp. 71–76, 1990.
- Y. Stylianou, O. Cappe, E. Moulines. Continuous probabilistic transform for voice conversion. IEEE Trans. Speech & Audio Process., Vol. 6, No. 2, pp. 131–142, 1998.
- L. Sun, K. Li, H. Wang, S. Kang, H.M. Meng. Phonetic posteriorgrams for many-to-one voice conversion without parallel data training. Proc. IEEE ICME, 6 pages, 2016.
- T. Toda, A.W. Black, K. Tokuda. Voice conversion based on maximum likelihood estimation of spectral parameter trajectory. IEEE Trans. Audio, Speech & Lang. Process., Vol. 15, No. 8, pp. 2222–2235, 2007.
- [http://spcc.csd.uoc.gr/SPCC2019/Lectures/SPCC2019\\_VC\\_Lecture\\_TomokiTODA.pdf](http://spcc.csd.uoc.gr/SPCC2019/Lectures/SPCC2019_VC_Lecture_TomokiTODA.pdf)
- Pantazis et al. , ” <https://www.csd.uoc.gr/~spcc/> ”



# Speech Intelligibility Enhancement

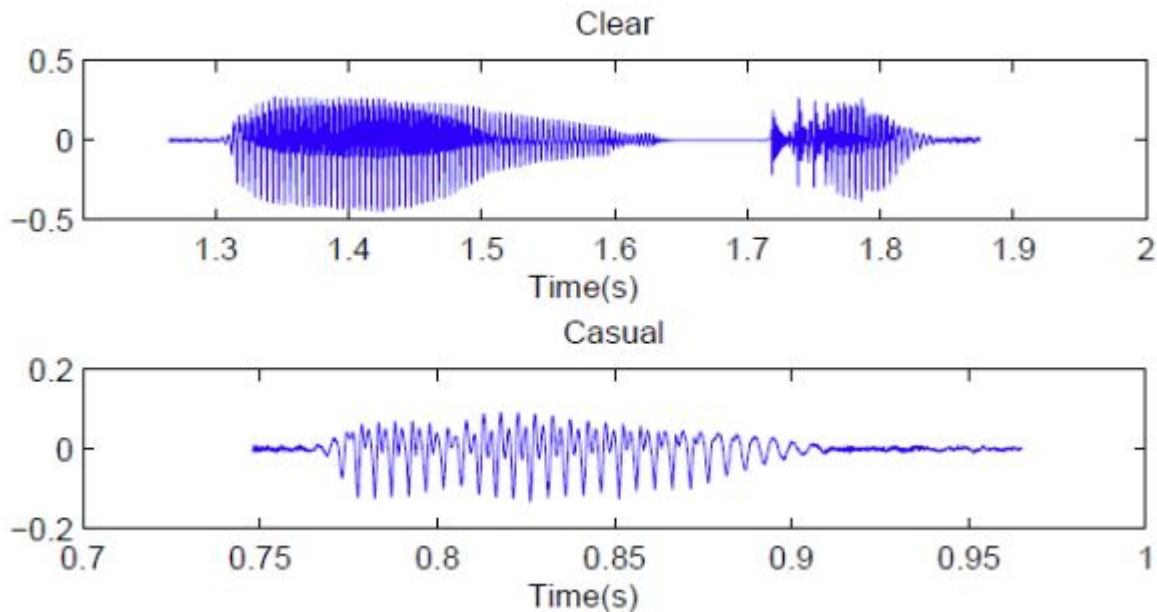
## From Casual to Clear Speech

# Clear and Conversational speech

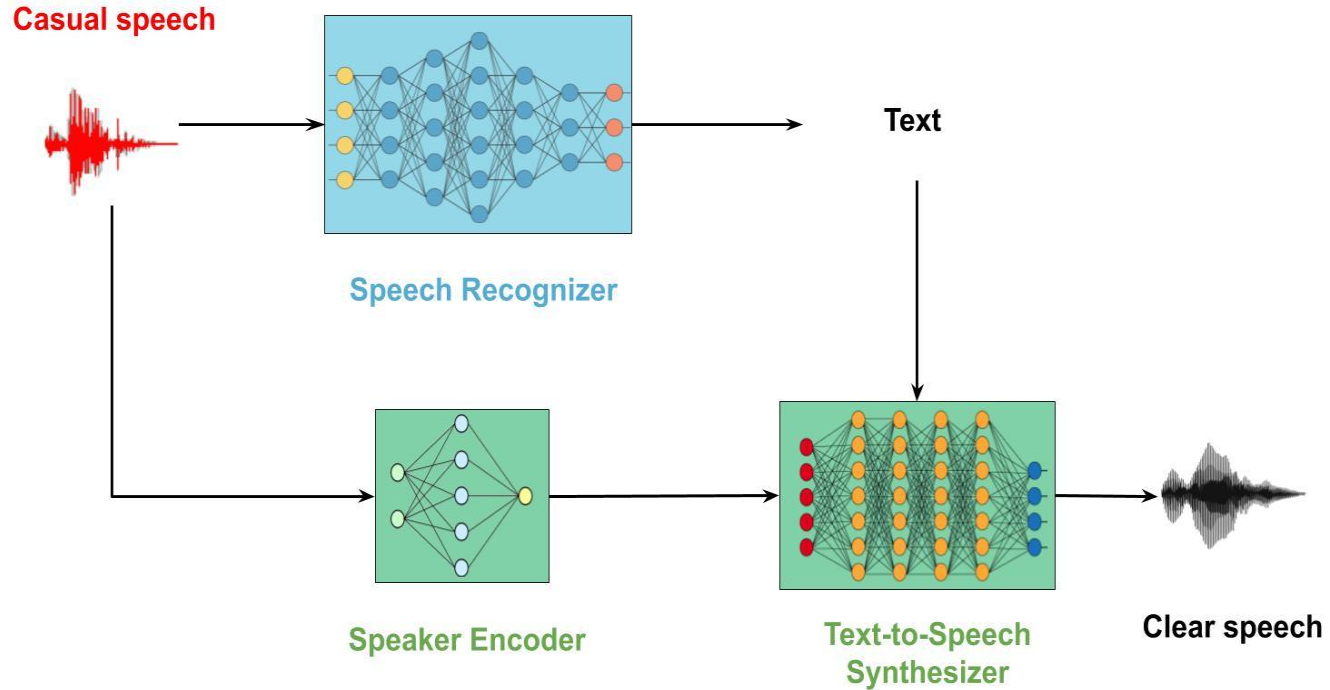
- Clear speech is a speaking style adopted by speakers in an attempt to maximize the clarity of their speech.
- Conversational speech is produced under casual or typical circumstances when no special speaking effort is made.
- However, in the presence of a communication difficulty, humans adopt different speaking styles.

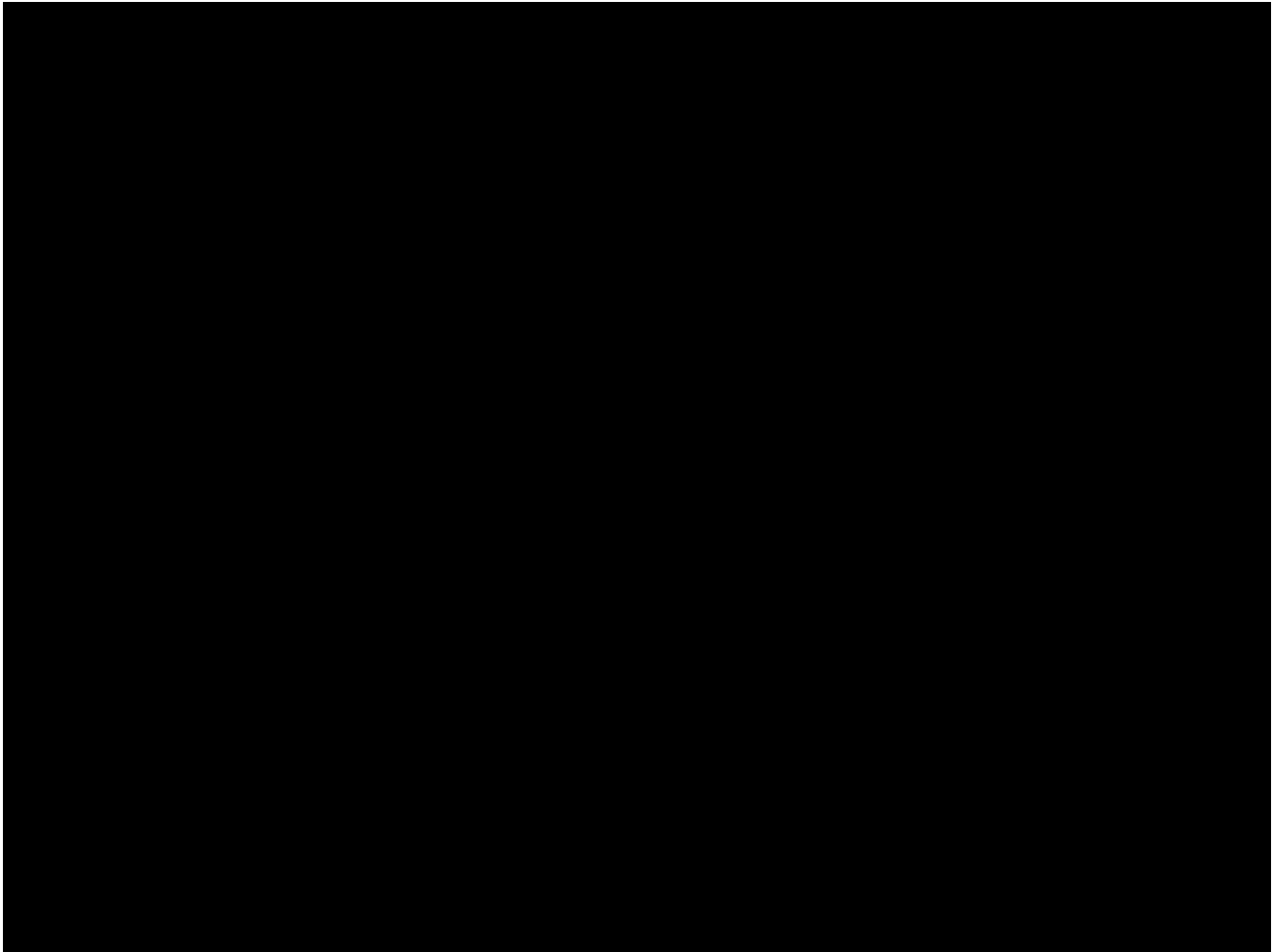
# Clear and Conversational speech

- The speaking style they adopt depends mostly on the communication barrier they want to overcome in order to communicate.



# System Design





# Sound Samples

<http://ixion.csd.uoc.gr/shifaspv/listest/index.php?n=Main.lcassp-show-tell>