

# DEEP LEARNING FOR SPEECH AND LANGUAGE

Winter School at UPC TelecomBCN Barcelona. 24-30 January 2018.



## Instructors



Marta R.  
Costa-jussà



José A. R.  
Fonollosa



Santiago  
Pascual



Javier  
Hernando



Antonio  
Bonafonte



Xavier  
Giró-i-Nieto

Organized by



UNIVERSITAT POLITÈCNICA  
DE CATALUNYA  
BARCELONATECH



Supported by



GitHub Education

Google Cloud Platform

+ info: <https://telecombcn-dl.github.io/2018-dsl/>

[\[course site\]](#)



#DLUPC

Day 4 Lecture 2

## Speech to speech paradigms



Santiago Pascual  
[santi.pascual@upc.edu](mailto:santi.pascual@upc.edu)

PhD Candidate  
Universitat Politècnica de Catalunya  
Technical University of Catalonia



# Outline

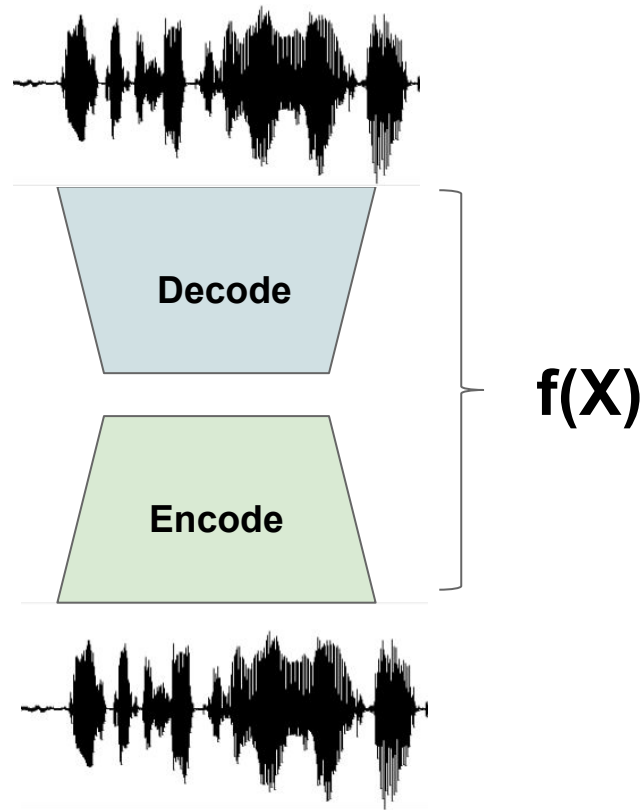
1. Introduction
2. Encoder-Decoder Paradigms
  - a. Generative modeling
3. Speech Enhancement
  - a. Discriminative Procedure
  - b. SEGAN/FSEGAN
4. Voice Conversion

# Introduction

# Speech to speech

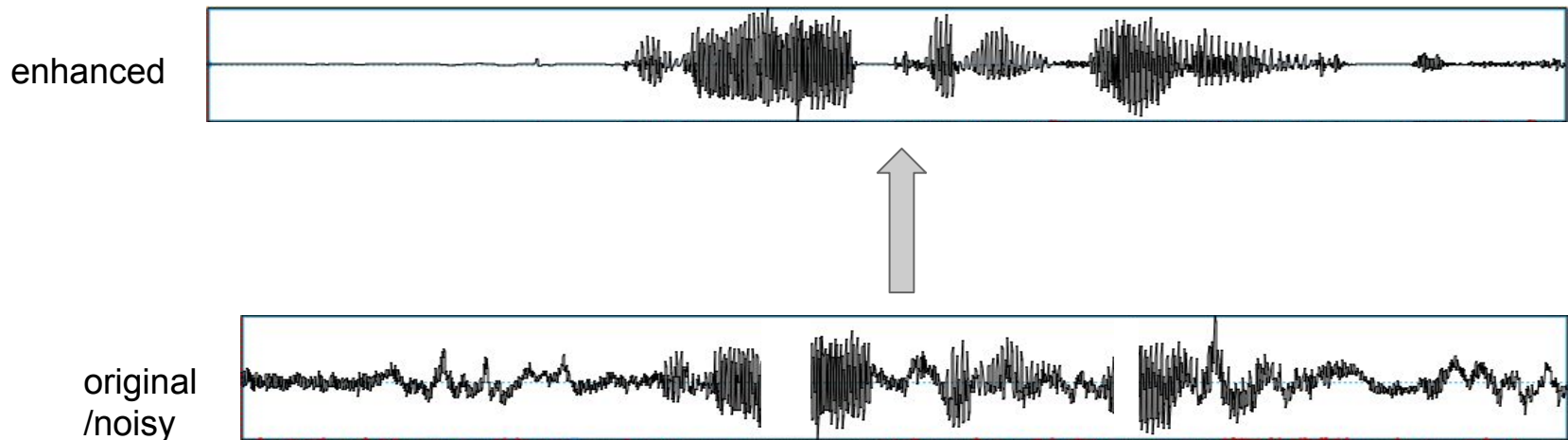
Speech is transformed through a non-linear function  $Y = f(X)$ :

- Enhance/Denoise signal
- Convert content respecting identity
  - Translation
- Convert identity respecting content
  - Voice Conversion



# Speech Enhancement/Denoising

Recover lost information or add enhancing details by learning the natural distribution of audio samples.



# Speech Enhancement

Applicable to many scenarios:

- Improving automatic speech recognition (ASR).
- Improve intelligibility in complex communication scenarios (like airplanes).
- For hearing aid implants.
- Enhance low quality recordings in speech synthesis data to train a system.

# Voice Conversion

Transfer the spoken contents and style from one speaker A to another speaker B.



Speaker A

"I am so happy"



**Voice Conversion**



"I am so happy"



Speaker B

# Voice Conversion

Also: transfer the spoken contents and style from within same speaker identity.



Speaker A

“We won...”



**Voice Conversion**

“We won!”



Speaker A



# Voice Conversion

## Potential Applications:

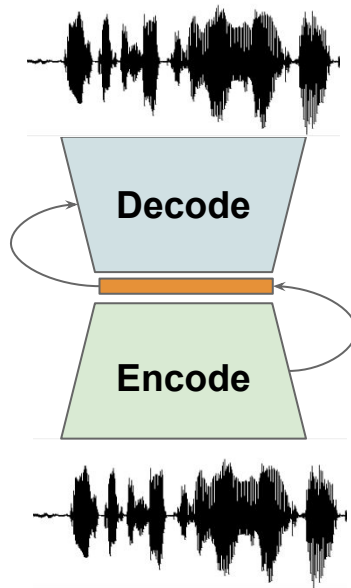
- Technologies to help people with motor speech disorders like dysarthria.
- Additional flexible block to speech synthesis systems, specially to unit selection ones, where we can enforce emotions and prosody changes.
- Dubbing industry. Human speech contains a set of expressive and natural patterns that are hard to obtain directly from text like in TTS.

# Encoder-Decoder Paradigms

# Encoder-Decoder paradigm

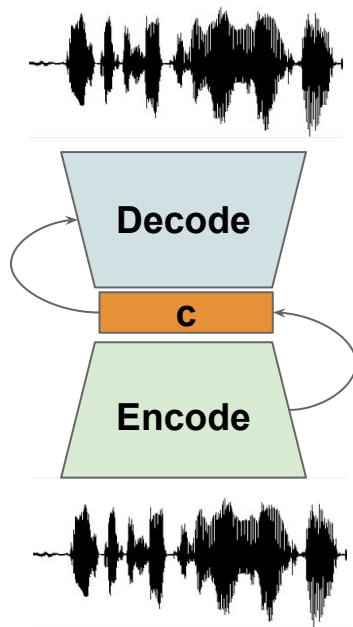
These speech2speech systems typically work under an encoder-decoder framework:

- Build an intermediate representation that captures latent characteristics of the spoken utterance.
- Reconstruct the signal with the proper new features.



# Vanilla AutoEncoders

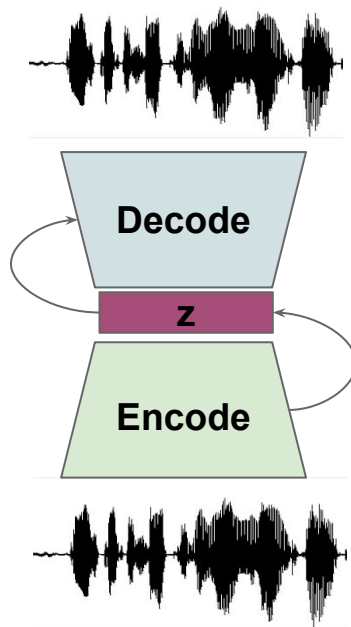
- Encoder mapping  $\mathbf{c} = E(\mathbf{x})$  is deterministic, as well as code vector  $\mathbf{c}$ .
- Decoder mapping reconstructs  $\mathbf{x}$  into a plausible version  $\hat{\mathbf{x}}$  deterministically.



# Variational AutoEncoders

([Kingma and Welling, 2014](#))

- Encoder mapping  $\mathbf{z} = E(\mathbf{x})$  is deterministic, but we apply restrictions on  $\mathbf{Z}$  space, so that it follows a prior probability density, like isotropic Normal one:  $N(0, I)$ .
- Decoder mapping reconstructs a sampled  $\mathbf{z}$  into a plausible version  $\mathbf{x}^\wedge$  deterministically.



NOTE: Working directly with waveforms is a very recent thing (1 year at most), and one of the most challenging parts of deep speech2speech systems.

# VQ-VAE

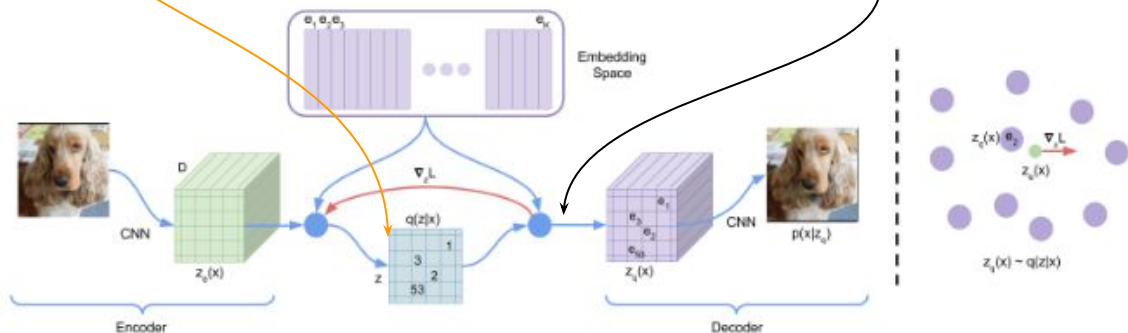
([Van den Oord et al. 2017](#))

- Z** space is a discretized embedding space, so every encoded point  $z(\mathbf{x})$  is mapped to nearest embedding  $\mathbf{e}$ , which is the information given to decode the sample.

$$q(z = k|x) = \begin{cases} 1 & \text{for } k = \operatorname{argmin}_j \|z_e(x) - e_j\|_2, \\ 0 & \text{otherwise} \end{cases},$$

$$z_q(x) = e_k, \quad \text{where } k = \operatorname{argmin}_j \|z_e(x) - e_j\|_2$$

$$e \in \mathbb{R}^{K \times D}$$



# Generative Adversarial Networks

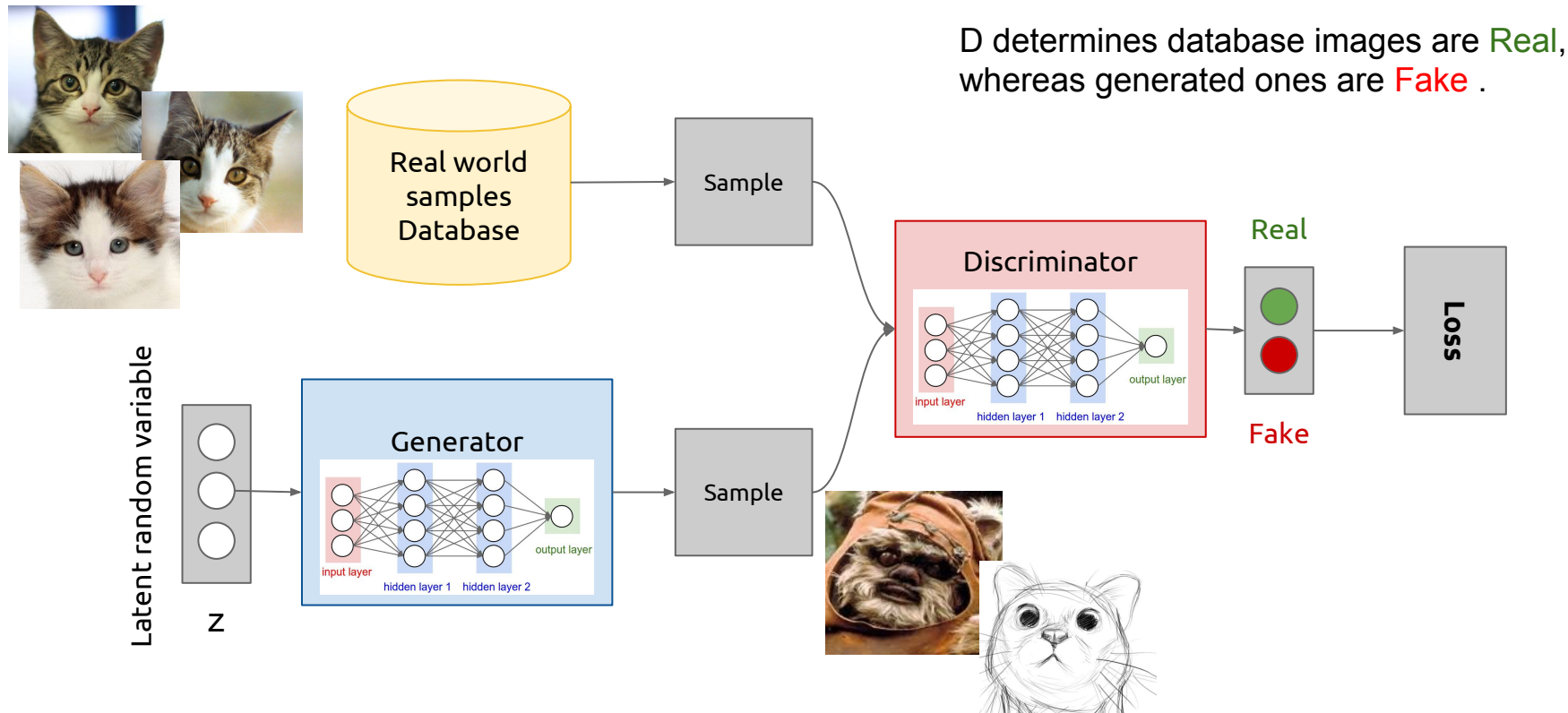
([Goodfellow et al. 2014](#))

We have two modules: **Generator** (G) and **Discriminator** (D).

- They “fight” against each other during training → **Adversarial Training**
- G mission: Fool D to misclassify.
- D mission: Discriminate between G samples and real samples.



# Generative Adversarial Networks



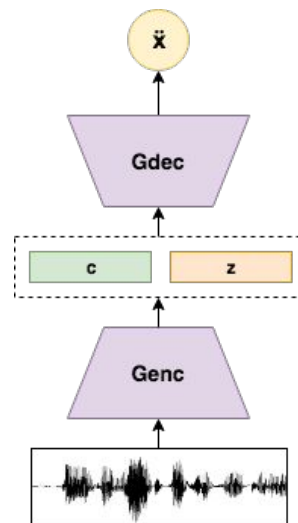
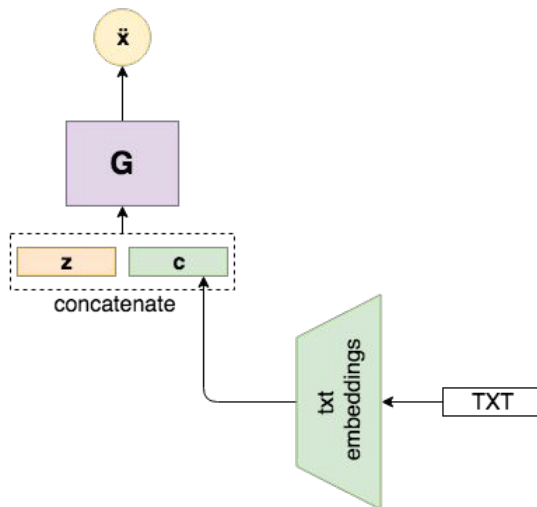
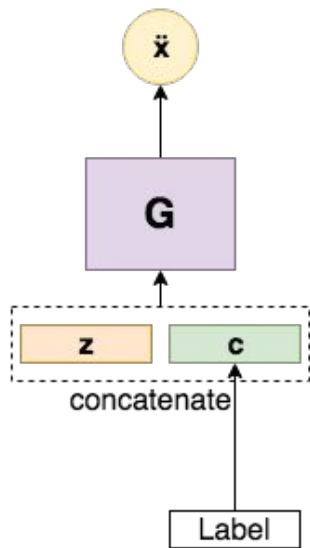


# Conditional GANs

For details on ways to condition GANs:  
[Ways of Conditioning Generative Adversarial Networks \(Wack et al.\)](#)

GANs can be conditioned on other info extra to **z**: text, labels, speech, etc..

**z** might capture random characteristics of the data (variabilities of plausible futures), whilst **c** would condition the deterministic parts !

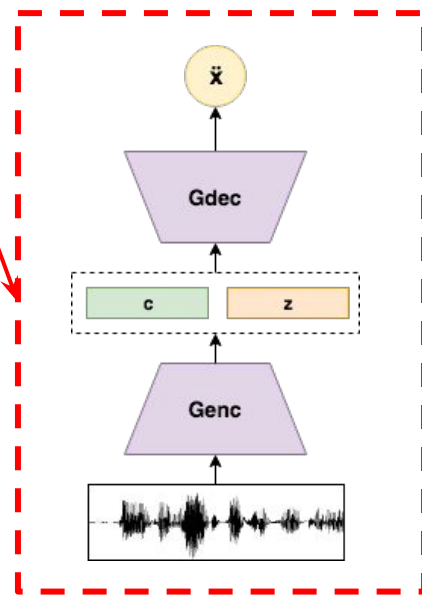
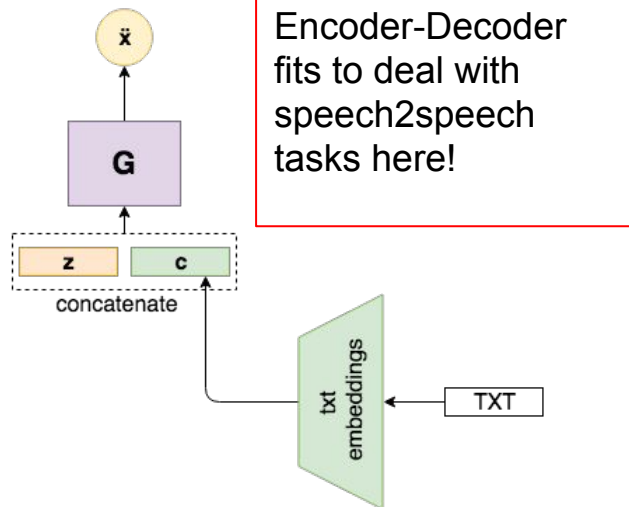
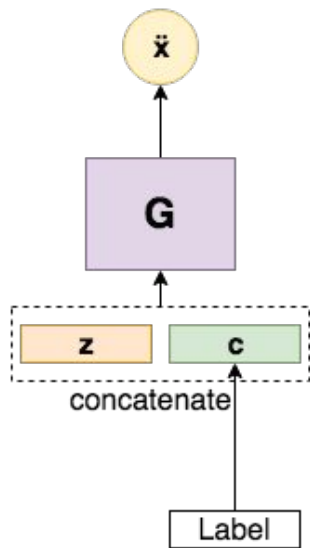


# Conditional GANs

For details on ways to condition GANs:  
[Ways of Conditioning Generative Adversarial Networks \(Wack et al.\)](#)

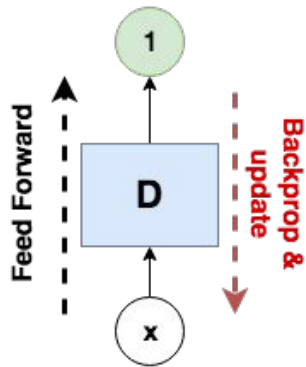
GANs can be conditioned on other info extra to  $\mathbf{z}$ : text, labels, speech, etc..

$\mathbf{z}$  might capture random characteristics of the data (variabilities of plausible futures), whilst  $\mathbf{c}$  would condition the deterministic parts !



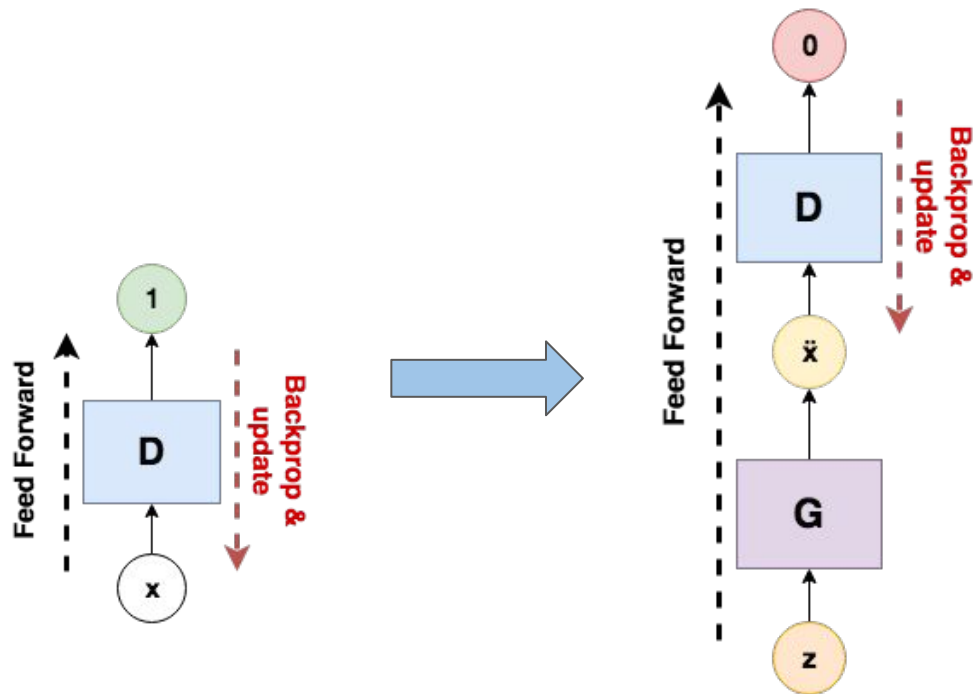
# Adversarial Training (1)

- Pick a sample  $\mathbf{x}$  from training set
- Show  $\mathbf{x}$  to **D** and update weights to output 1 (real)



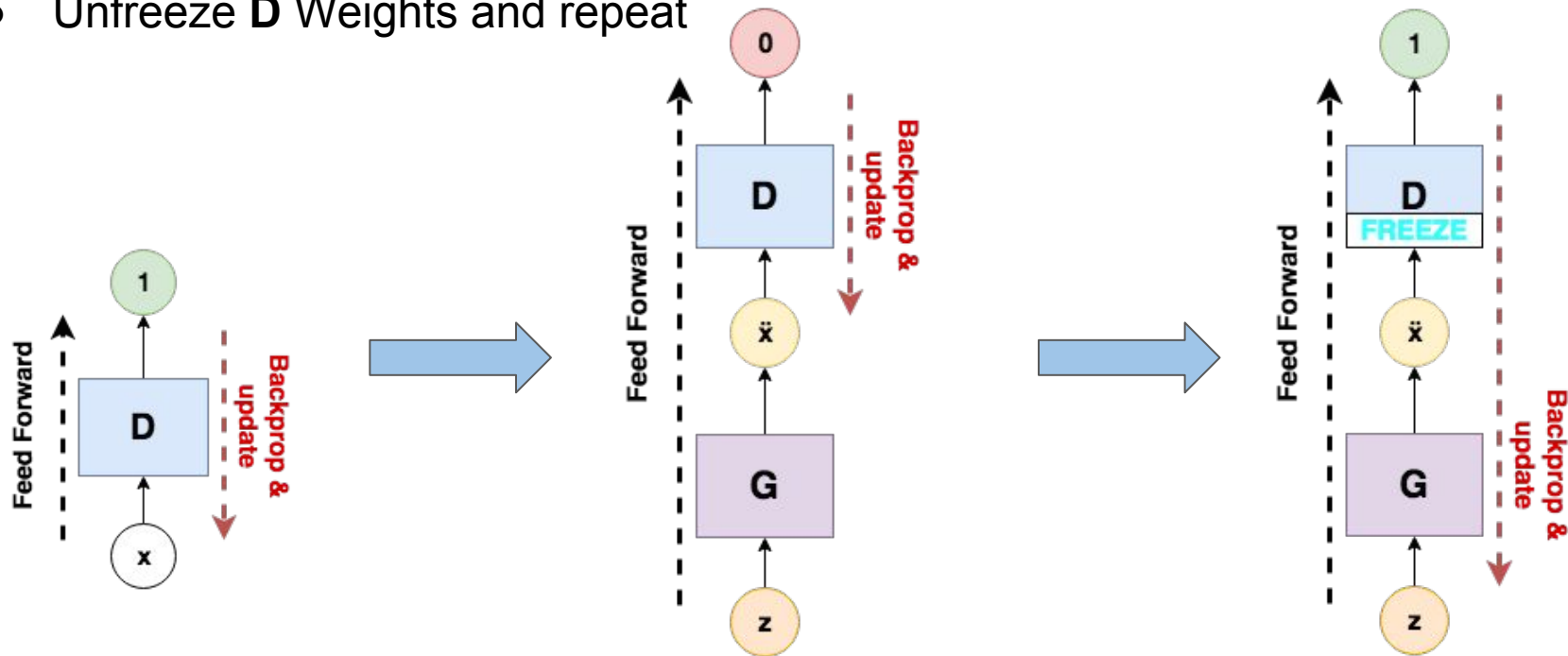
# Adversarial Training (2)

- $G$  maps sample  $z$  to  $\tilde{x}$
- show  $\tilde{x}$  and update weights to output 0 (fake)



# Adversarial Training (3)

- Freeze **D** weights
- Update **G** weights to make **D** output 1 (just **G** weights!)
- Unfreeze **D** Weights and repeat



# Least Squares GAN

Main idea: shift to loss function that provides smooth & non-saturating gradients in D

- **Because of sigmoid saturation** in binary classification loss, G gets no info when D gets to label true examples → **vanishing gradients make G no learn**
- Least squares loss improves learning with notion of distance of  $P_{model}$  to  $P_{data}$ :

$$\begin{aligned}\min_D V_{LSGAN}(D) &= \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})} [(D(\mathbf{x}) - 1)^2] + \frac{1}{2} \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [(D(G(\mathbf{z})))^2] \\ \min_G V_{LSGAN}(G) &= \frac{1}{2} \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [(D(G(\mathbf{z})) - 1)^2],\end{aligned}$$

# Voice Conversion

# Parallel corpora and frame-wise VC

General VC pipeline with Discriminative model:

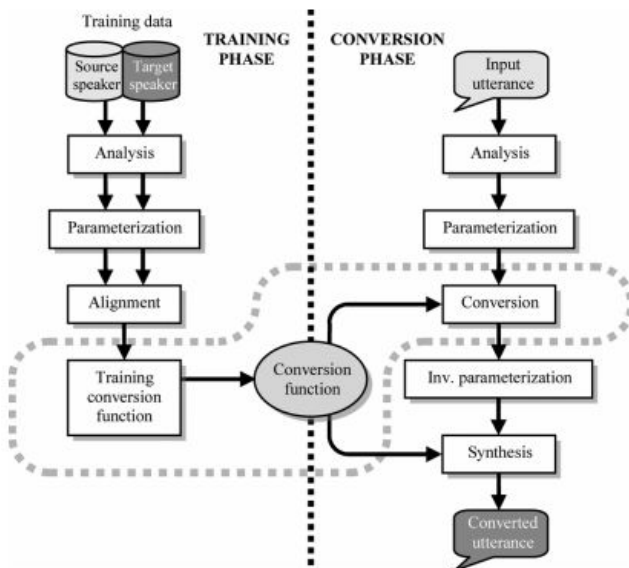


Figure credit: Daniel Erro



# Parallel corpora and frame-wise VC

General VC pipeline with Discriminative model:

(1) Spectral features are extracted

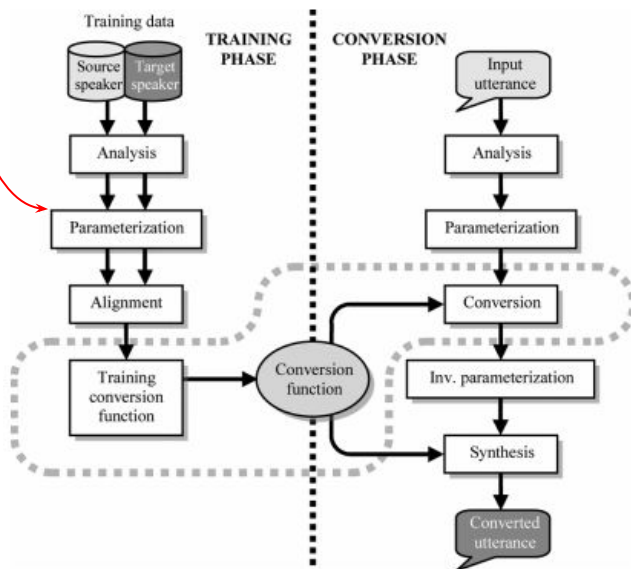
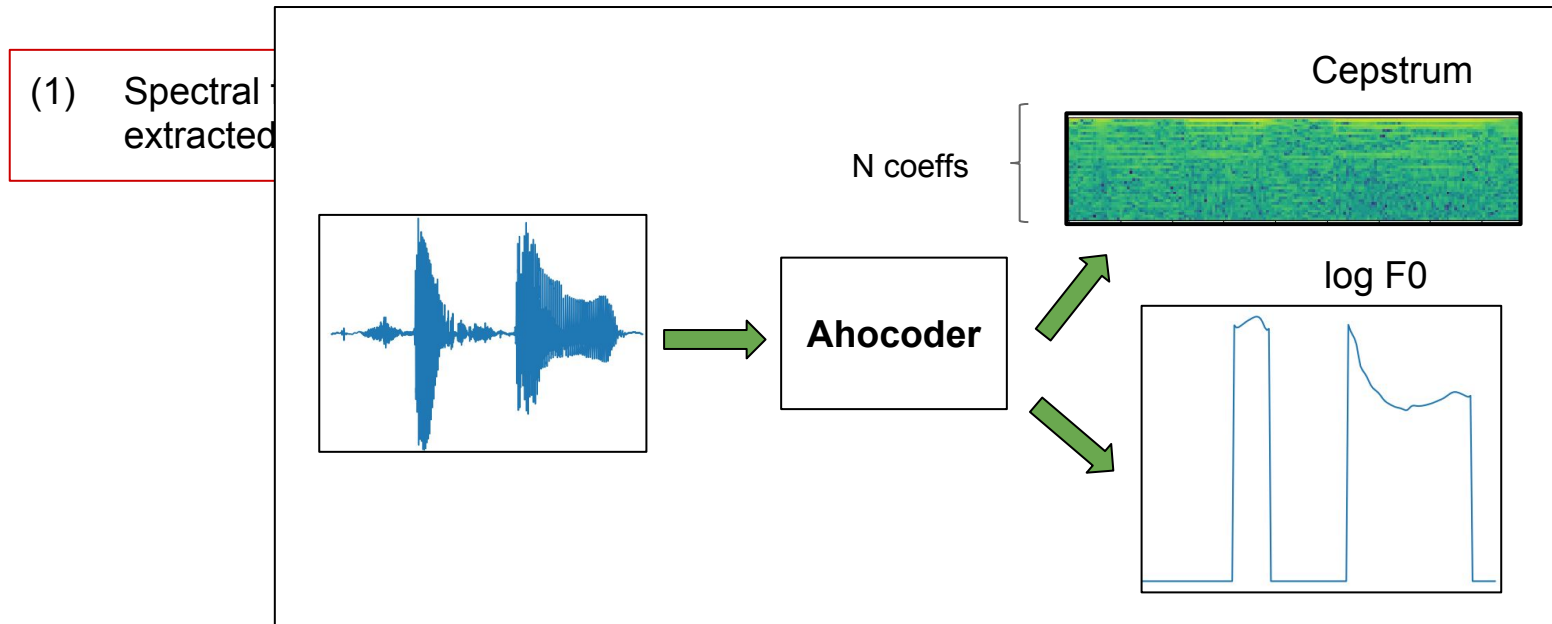


Figure credit: Daniel Erro

# Parallel corpora and frame-wise VC

General VC pipeline with Discriminative model:



TRAIN

# Parallel corpora and frame-wise VC

General VC pipeline with Discriminative model:

(2) Alignment process in training data: Dynamic Time Warping

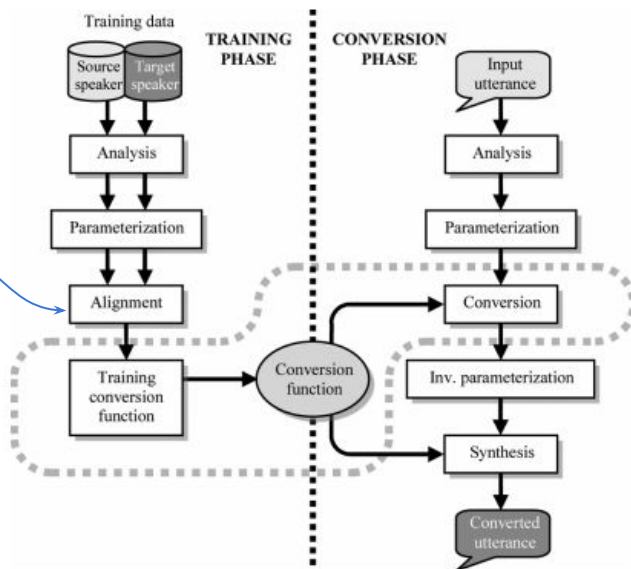
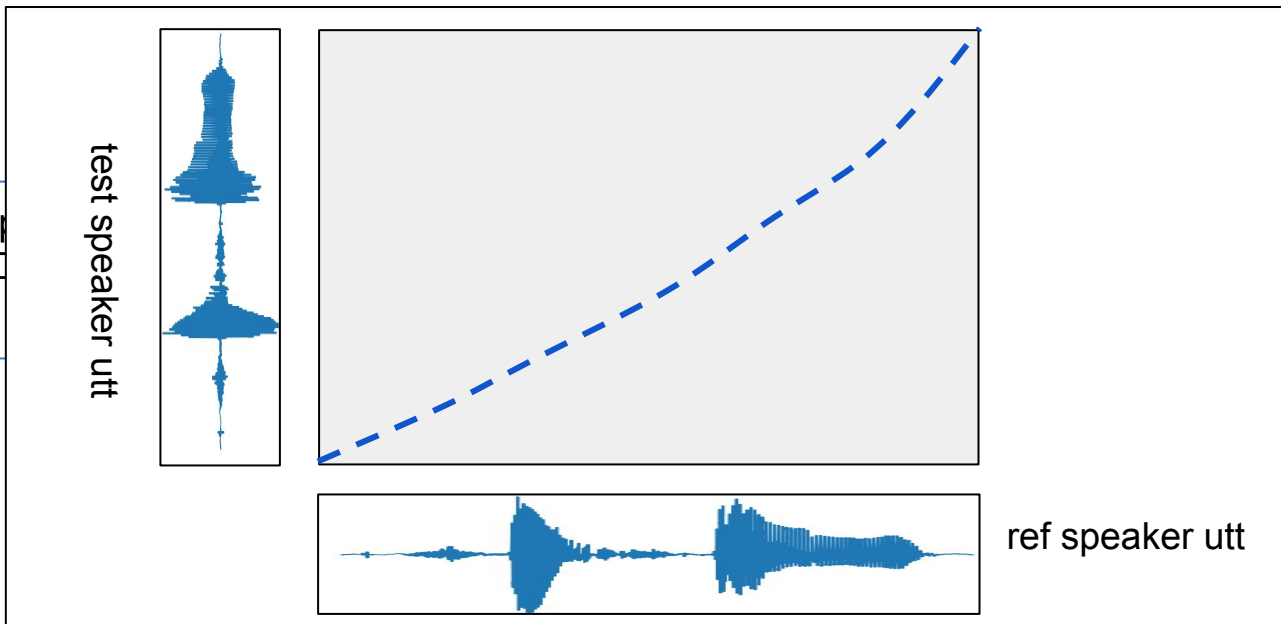


Figure credit: Daniel Erro

# Parallel corpora and frame-wise VC

General VC pipeline with Discriminative model:

(2) Alignment  
training data: D  
Time Warping



TRAIN

# Parallel corpora and frame-wise VC

General VC pipeline with Discriminative model:

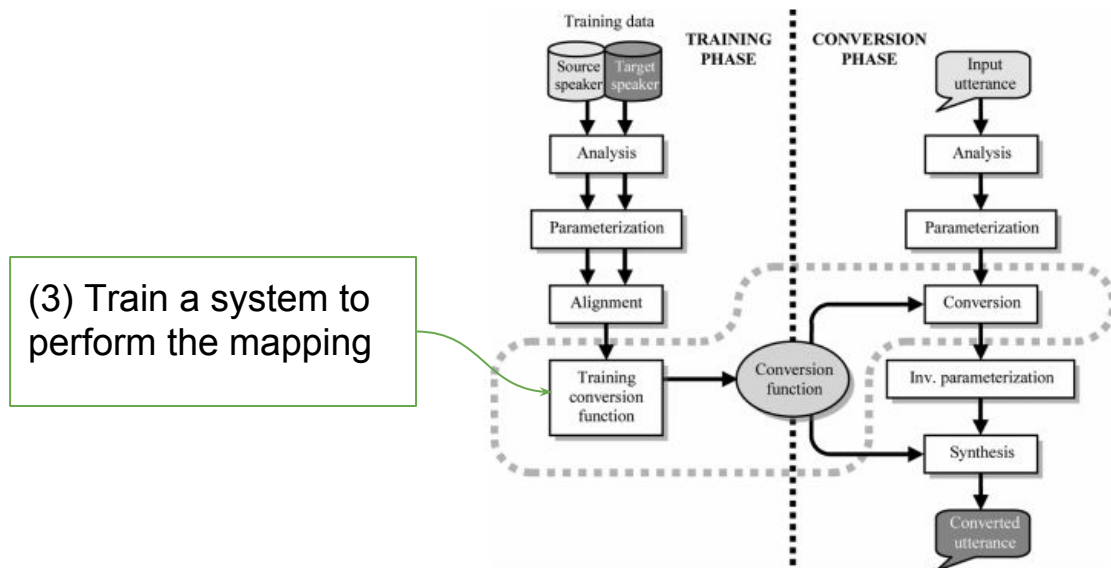


Figure credit: Daniel Erro

# Parallel corpora and frame-wise VC

General VC pipeline with Discriminative model:

Pitch can be linearly converted, pre-calculating both speakers' (source and target) statistical moments (mean and variance) among sliding window frames in training set:

$$\log(f0_{conv}) = \mu_{tgt} + \frac{\sigma_{tgt}}{\sigma_{src}}(\log(f0_{src}) - \mu_{src})$$

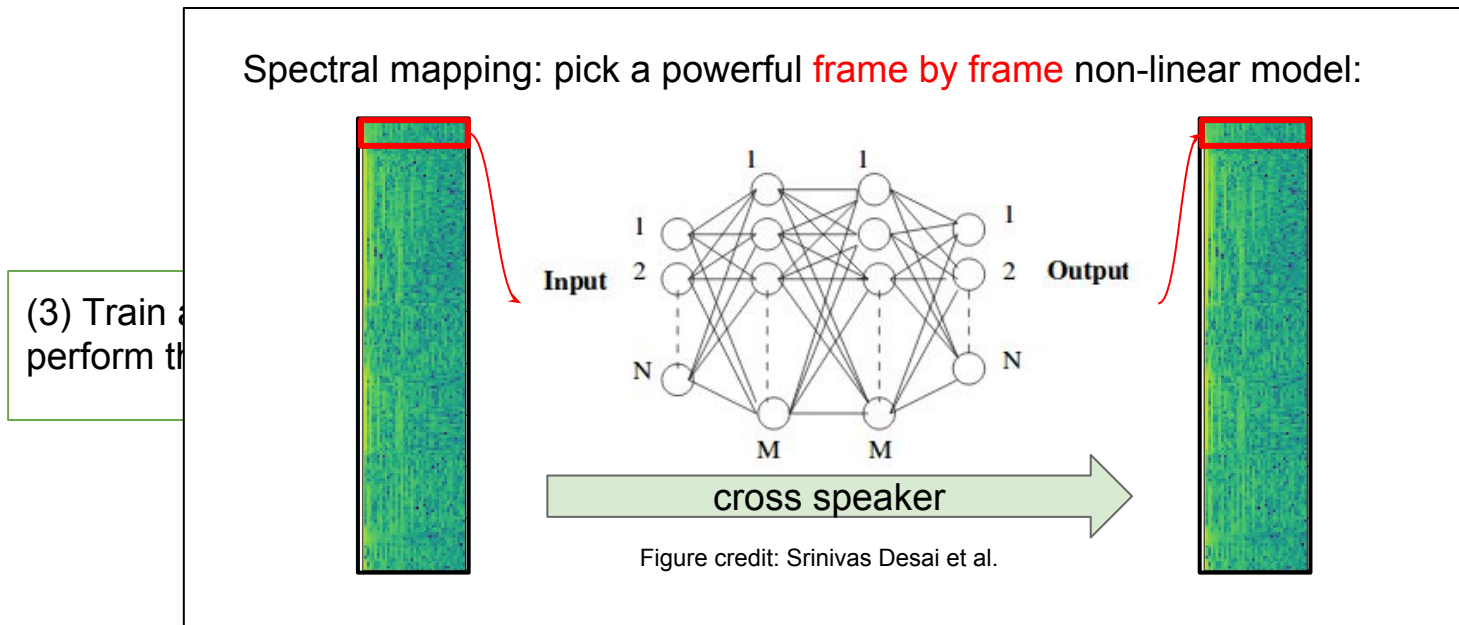
(3) Train &  
perform the

TRAIN

Figure credit: Daniel Erro

# Parallel corpora and frame-wise VC

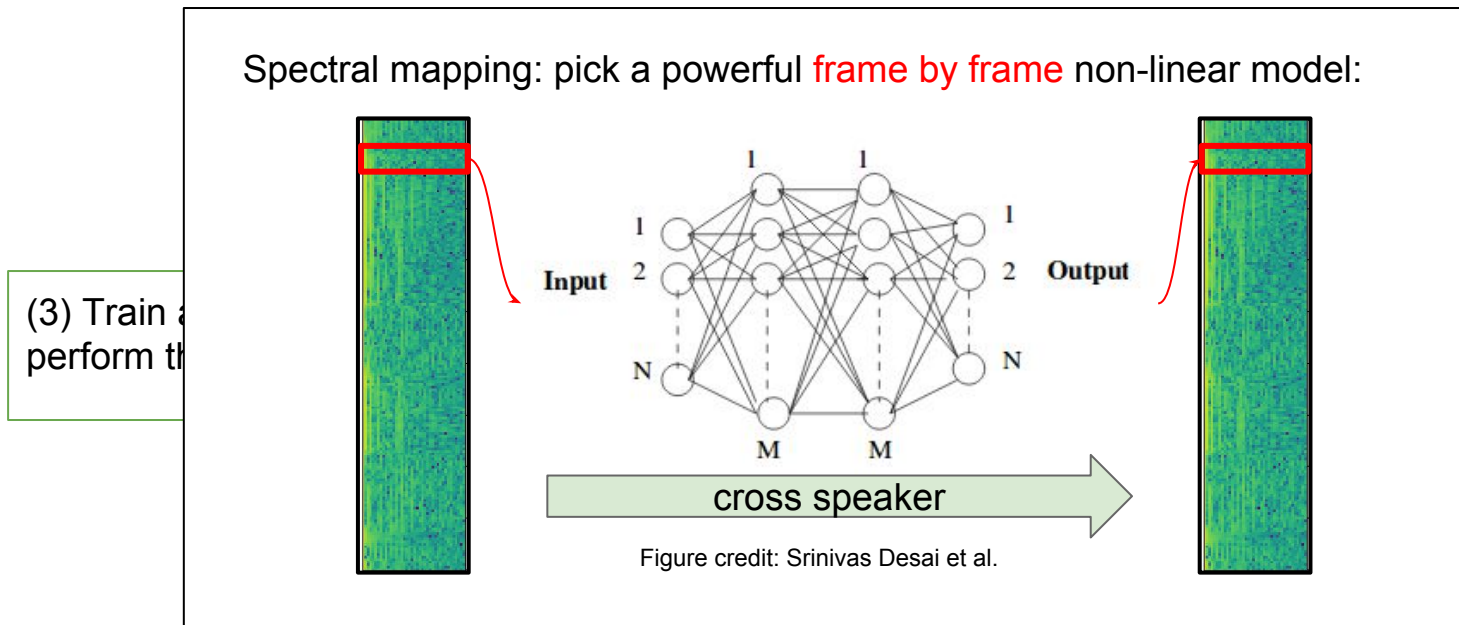
General VC pipeline with Discriminative model:



TRAIN

# Parallel corpora and frame-wise VC

General VC pipeline with Discriminative model:

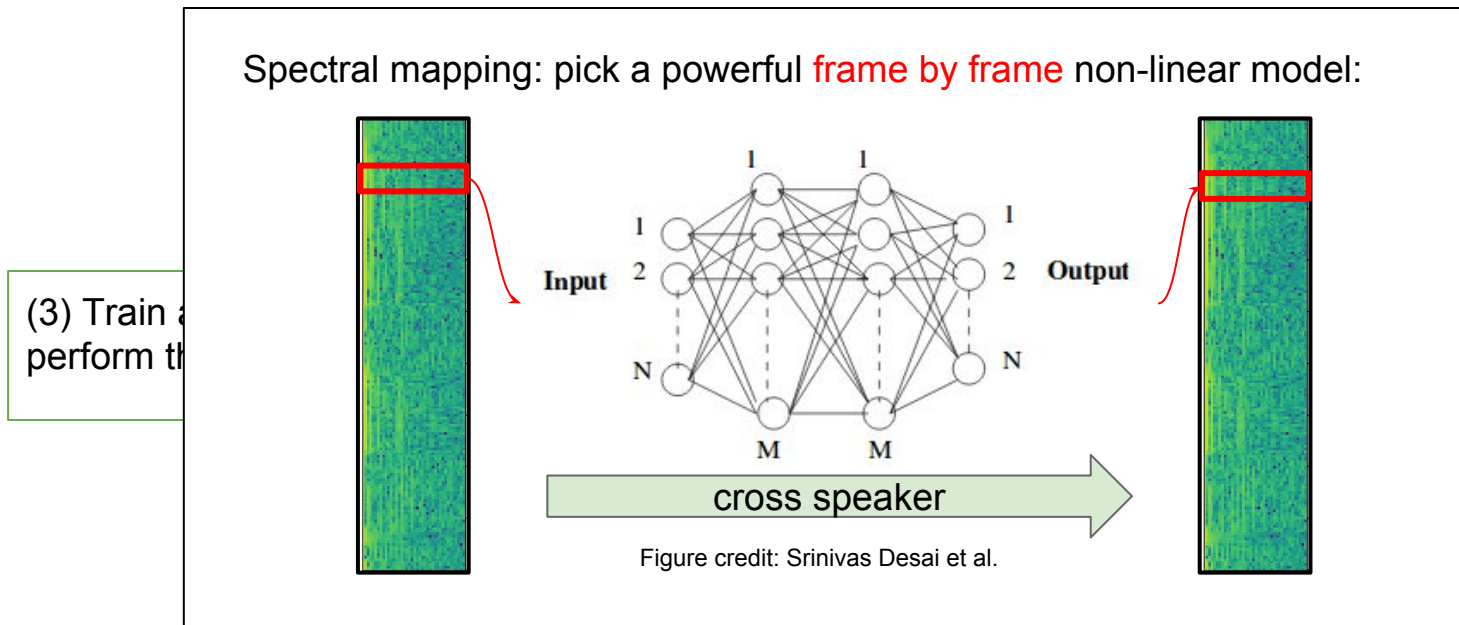


TRAIN



# Parallel corpora and frame-wise VC

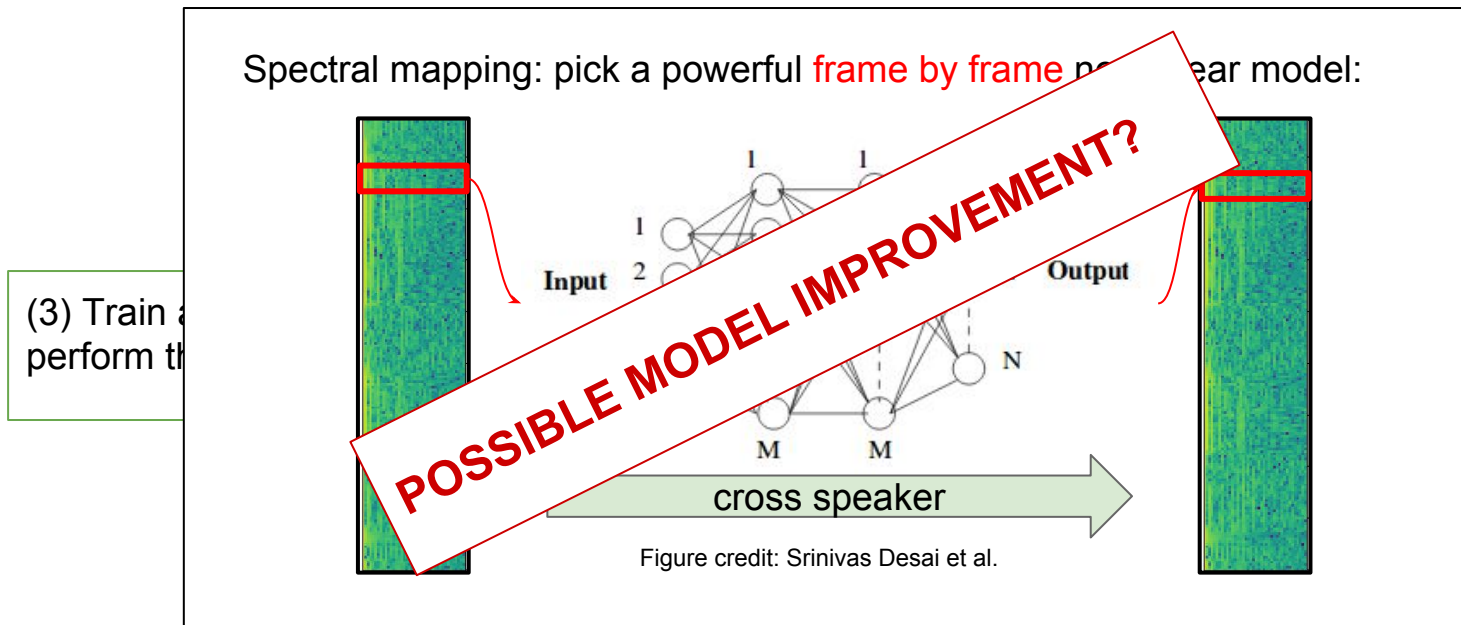
General VC pipeline with Discriminative model:



TRAIN

# Parallel corpora and frame-wise VC

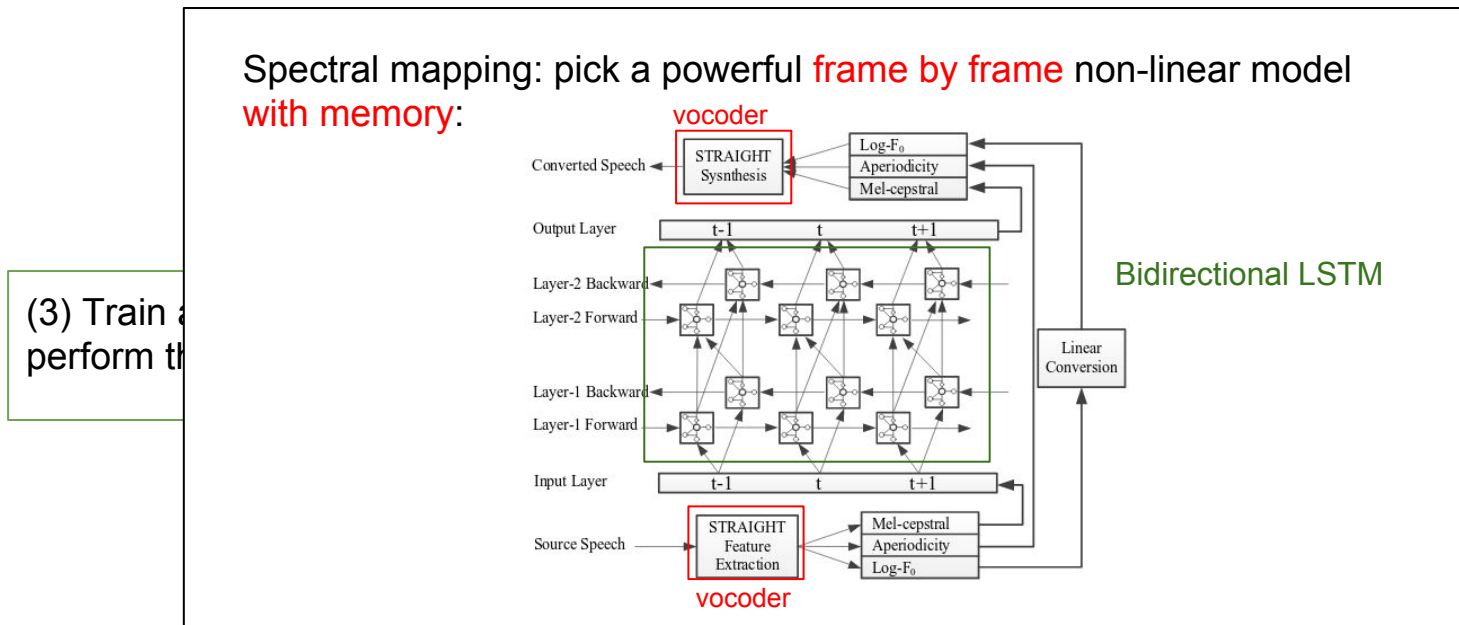
General VC pipeline with Discriminative model:



TRAIN

# Parallel corpora and frame-wise VC [\(Sun et al. 2015\)](#)

General VC pipeline with Discriminative model:



TRAIN

# Parallel corpora and frame-wise VC

General VC pipeline with Discriminative model:

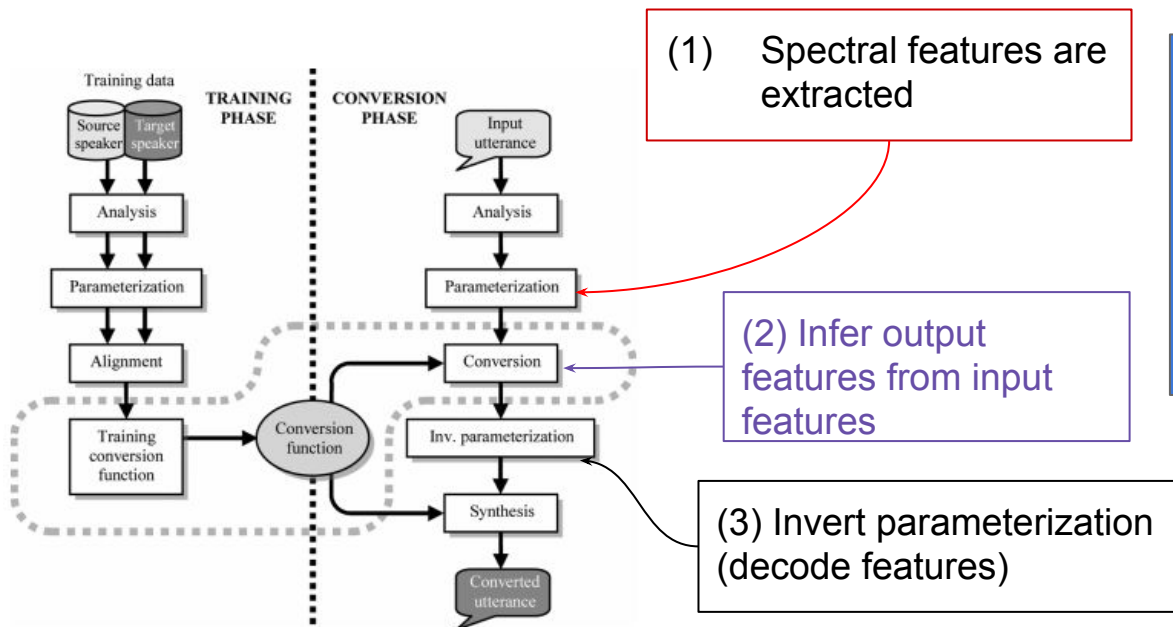
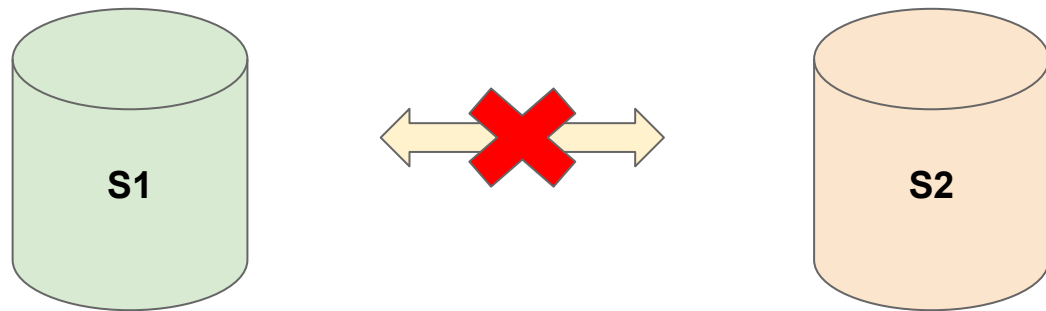


Figure credit: Daniel Erró

# Unaligned corpora

- Speakers do NOT say the same, so there's no content to align.
- Speakers can even speak in different languages!



Challenging transferrability problem: no supervised discriminative approach

# VAE based VC

([Hsu et al. 2016](#))

We can take advantage of Variational Auto-Encoder training procedure to learn latent representations of speakers, and a deterministic identity code will map all back to destination acoustic space.

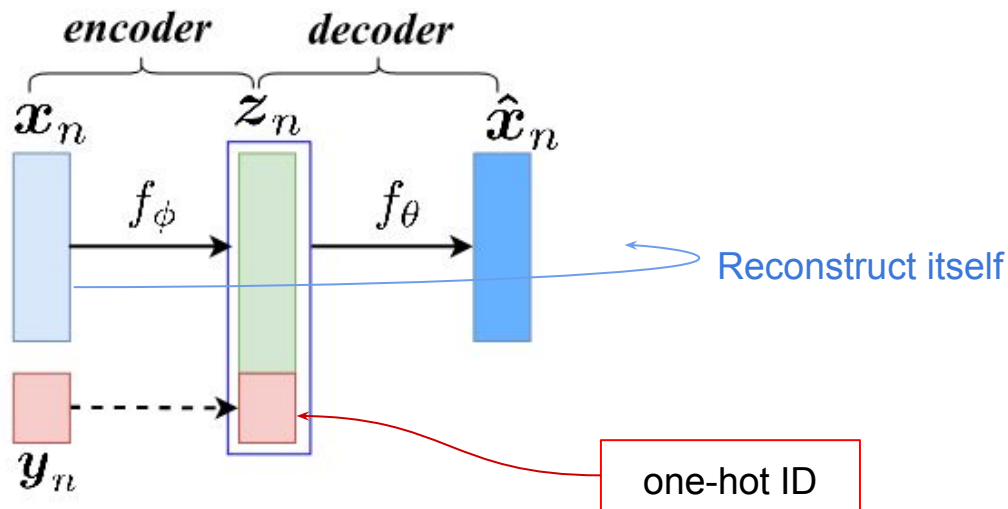


Figure credit: Hsu et al.

# Vector Quantised-VAE (end to end) [samples](#)

Latest most successful and natural sounding approach has been VQ-VAE by Google DeepMind. They build a discrete latent space that resembles a phoneset unsupervisedly! A **Wavenet** decodes the latent codes **conditioned on one-hot ID**.

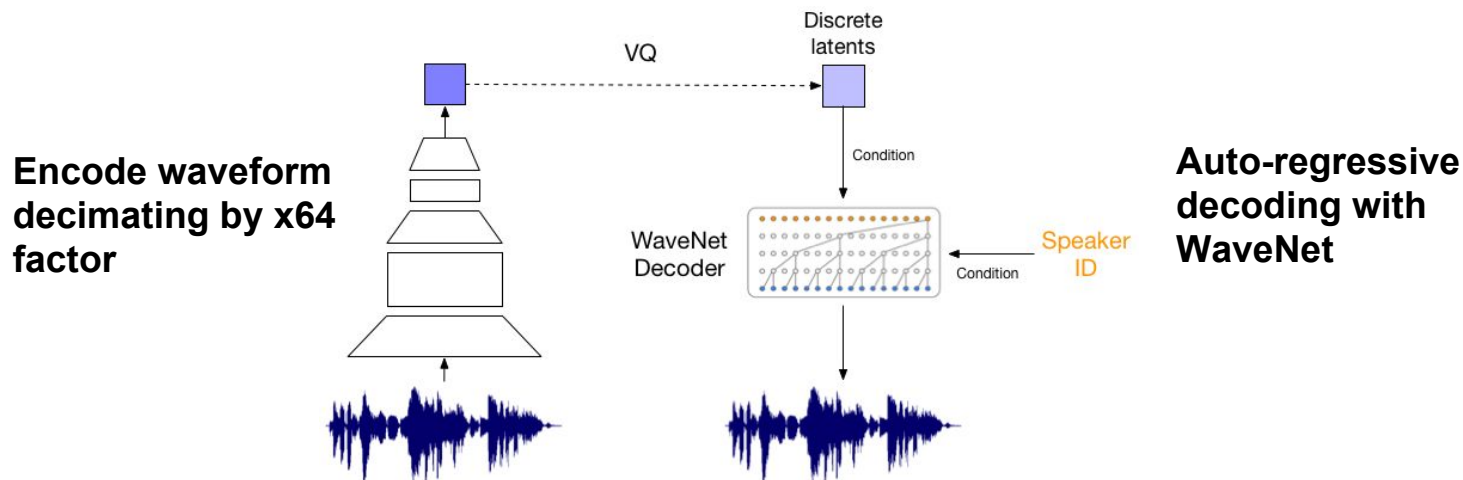


Figure credit: Aaron van den Oord

# VC Evaluation

- Typically subjective evaluation: like Mean Opinion Score (MOS) [1, 5] pooling a group of listeners opinions' in terms of (1) naturalness and (2) similarity to target.
- Objective metrics for specific features (e.g. Mel Cepstral Distortion [dB] for cepstrums, or RMSE [Hz] for F0 can serve as a guidance, but not as a final decision).



# Speech Enhancement

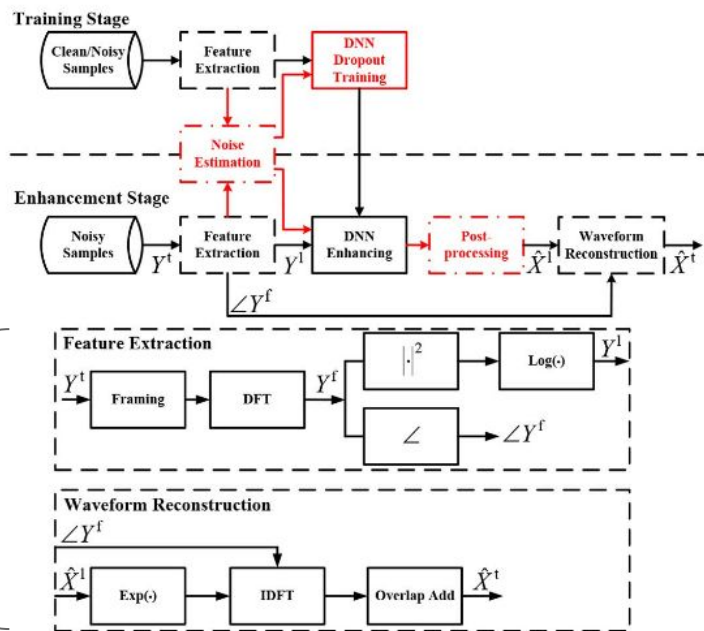
# SE Approaches

- Spectral subtraction: estimate noise activity during non-speech regions and subtract it.
- Subspace algorithms: decompose the higher dimensional noisy signal into a lower dimensional one where clean version lays.
- Spectral masking: predict a binary freq-time mask that can cancel out noisy bins.
- **Statistical model based: predict the clean features/signal as a statistical regression problem.**

# Discriminative regression

([Xu et al. 2015](#))

A DNN is used to map noisy parameterized speech (features) into the clean version as a regression problem (MSE estimation).



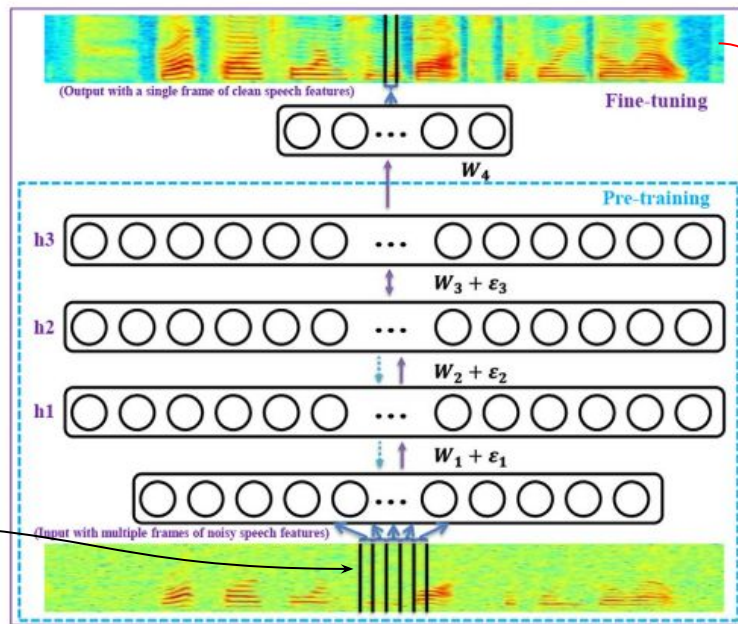
The log power of spectral module is enhanced (predicted). Phase remains the same and ISTFT recovers signal back.

Figure credit: Xu et al.

# Discriminative regression

([Xu et al. 2015](#))

A DNN is used to map noisy parameterized speech (features) into the clean version as a regression problem (**MSE** estimation).



$$Er = \frac{1}{N} \sum_{n=1}^N \|\hat{\mathbf{X}}_n(\mathbf{Y}_{n-\tau}^{n+\tau}, \mathbf{W}, \mathbf{b}) - \mathbf{X}_n\|_2^2.$$

Many input frames consider time correlation

Pre-training can be performed stacking RBMs as in paper's proposal.

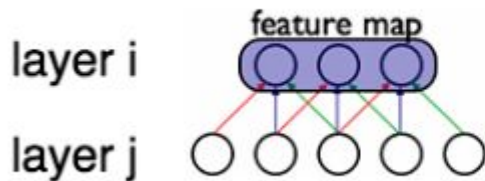
Two stages in Generator (fully convolutional) network:

1. Encoder (Downconv): Project noisy signal into a deterministic representation  $\mathbf{c}$  and concatenate to latent variable  $\mathbf{z} \sim N(0, \mathbf{I})$
2. Decoder (Deconv): Interpolate the intermediate hidden features w/ learnable params. until re-generation of clean speech.



# SEGAN: underlying structures

- 1D convolutional neural networks

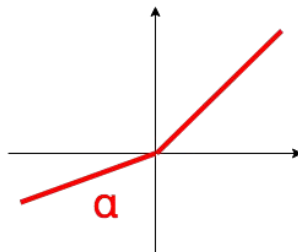


$$h_i^k = \tanh(W_{ij}^k * h_j + b_i^k)$$

$x$

- Virtual Batch Normalization: normalize layer responses with statistics from (reference\_batch + current\_batch) → less intra dependent statistics to avoid GAN instability.

- LeakyReLU/ParametricReLU:
  - $\alpha$  fixed (0.3) or learnable

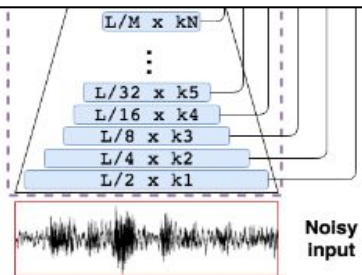


# SEGAN end to end training

- Show pairs of signals to “learn” a reconstruction loss.
- Use of L1 regularization to guide the GAN training.

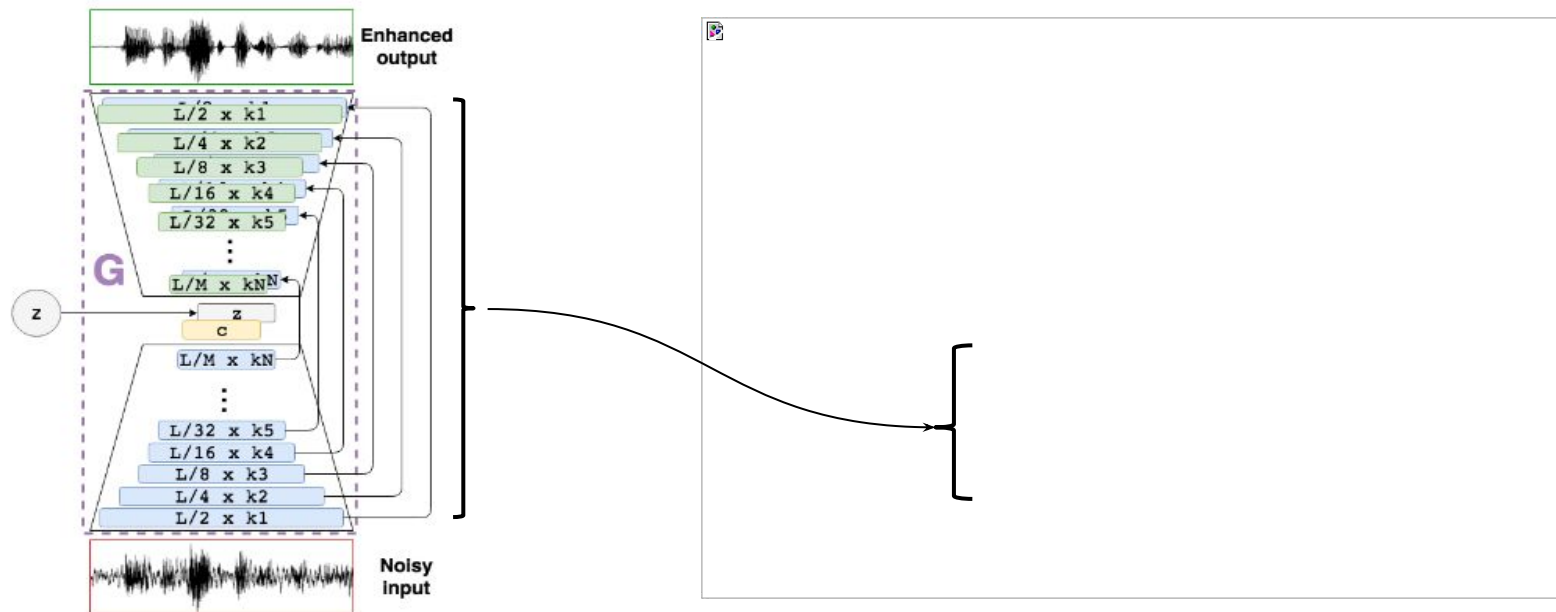
Final G loss: LSGAN  
Adversarial + weighted L1  
regularization/regression

$$\min_G V_{\text{LSGAN}}(G) = \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x}_c), \mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [(D(G(\mathbf{z}, \mathbf{x}_c)) - 1)^2] + \lambda \|G(\mathbf{z}, \tilde{\mathbf{x}}) - \mathbf{x}\|_1.$$



# SEGAN end to end training

- Show pairs of signals to “learn” a reconstruction loss.
- Use of L1 regularization to guide the GAN training.

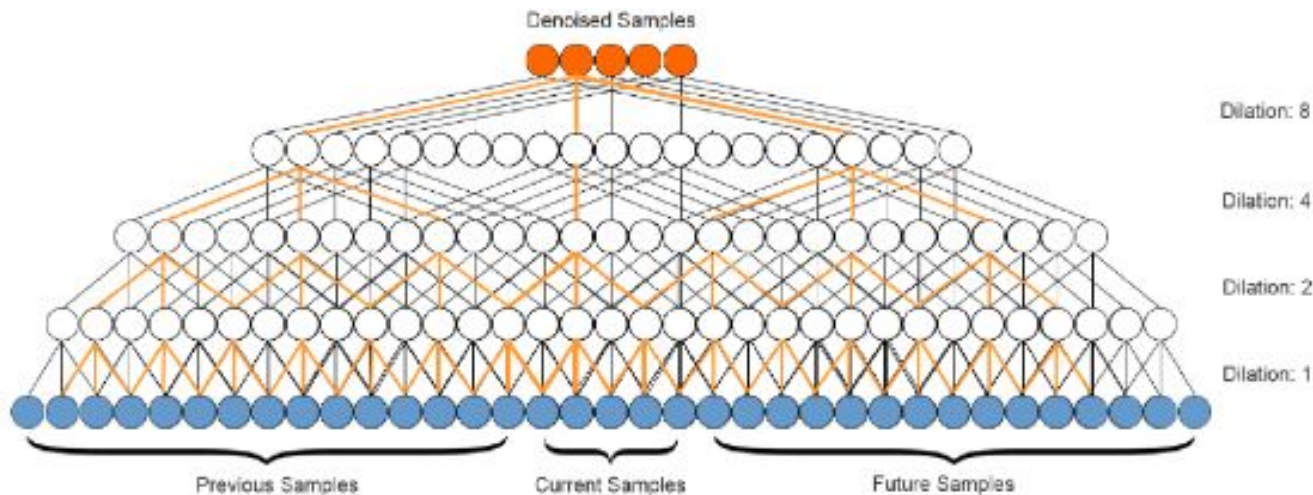




# Wavenet for Speech Denoising

([Rethage et al. 2017](#))

Wavenet proved to be effective as a generative model for raw speech and audio. A modified version of it was applied to speech denoising too, getting rid of the original autoregressive behavior, and dealing with a regression problem!



# Advanced SE research

Other active advances focus on using **perceptually weighted losses**, or using **enhancement as an internal stage** within another task, like Text-to-Speech (TTS) or Automatic Speech Recognition (ASR):

- [RNN-based SE for noise-robust TTS \(Valentini et al. 2016\)](#)
- [Perception Optimized Deep Denoising AutoEncoders for Speech Enhancement \(Gurunath and Georgiou 2016\)](#)
- [Exploring Speech Enhancement with Generative Adversarial Networks for Robust Speech Recognition \(Donahue et al. 2017\)](#)

# SE Evaluation

Typical objective metrics:

- PESQ: Perceptual Evaluation of Speech Quality [-0.5, 4.5]: designed for telephonic compression assessment.
- COVL: MOS prediction of the overall effect [1, 5]
  - CSIG: Mean opinion score (MOS) prediction of the signal distortion attending only to the speech signal [1, 5].
  - CBAK: MOS prediction of the intrusiveness of background noise [1, 5].
- SSNR: Segmental SNR [0, inf).

Nonetheless, subjective eval is always preferable (in any speech synthesis task)!

# Summary

- Speech2speech paradigms have been discussed, emphasizing the two salient ones at the moment: enhancement and conversion. All these methods are converging to end-to-end approaches.
- Voice Conversion parallel and non-parallel approaches have been reviewed, from classic frame-by-frame analysis to end-to-end VQ-VAE.
- Speech Enhancement methods have been reviewed, specially end-to-end ones, like SEGAN and Denoising Wavenet.
- Speech Enhancement is being included as an inherent end-to-end component for ASR and TTS, among others.
- Speech2speech paradigms are gaining momentum, specially the end-to-end embedded versions to process speech signals in real time in our handset devices.

# References

- [Auto-Encoding Variational Bayes \(Kingma and Welling 2014\)](#)
- [Generative Adversarial Networks \(Goodfellow et al. 2014\)](#)
- [Voice Conversion Using Artificial Neural Networks \(Desai et al. 2009\)](#)
- [Voice Conversion Using Deep Bidirectional Long-Short Term Memory Based Recurrent Neural Networks \(Sun et al. 2015\)](#)
- [Voice Conversion from Non-Parallel Corpora Using Variational Auto-encoder \(Hsu et al. 2016\)](#)
- [Neural Discrete Representation Learning \(van den Oord et al. 2017\)](#)
- [A Regression approach to speech enhancement based on deep neural networks \(Xu et al. 2015\)](#)
- [Perception Optimized Deep Denoising AutoEncoders for Speech Enhancement \(Gurunath and Georgiou 2016\)](#)
- [RNN-based SE for noise-robust TTS \(Valentini et al. 2016\)](#)
- [SEGAN: Speech Enhancement Generative Adversarial Network \(Pascual et al. 2017\)](#)
- [Exploring Speech Enhancement with Generative Adversarial Networks for Robust Speech Recognition \(Donahue et al. 2017\)](#)
- [A Wavenet for Speech Denoising \(Rethage et al. 2017\)](#)