

DEEP LEARNING FOR SPEECH & LANGUAGE

Winter Seminar UPC TelecomBCN, 24 - 31 January 2017

Instructors



Antonio Bonafonte J. Adrián Rodríguez Fonollosa Marta R. Costa-jussà Javier Hernando Santiago Pascual Elisa Sayrol Xavier Giró

Organizers



Image Processing Group
Signal Theory and Communications Department



UNIVERSITAT POLITÈCNICA
DE CATALUNYA
BARCELONATECH

+ info: [TelecomBCN.DeepLearning.Barcelona](https://www.telecombcn.com/deeplearning-barcelona)

[\[course site\]](#)

Day 3 Lecture 5

Parametric Speech Synthesis

Antonio Bonafonte



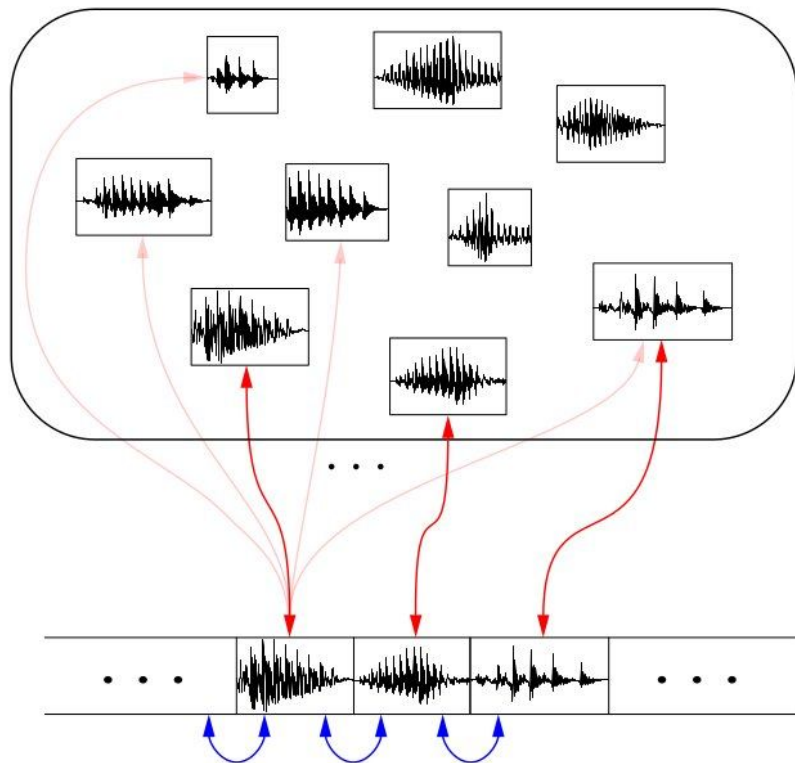
Main TTS Technologies

Concatenative speech synthesis + Unit Selection

Concatenate *best* prerecorded speech *units*

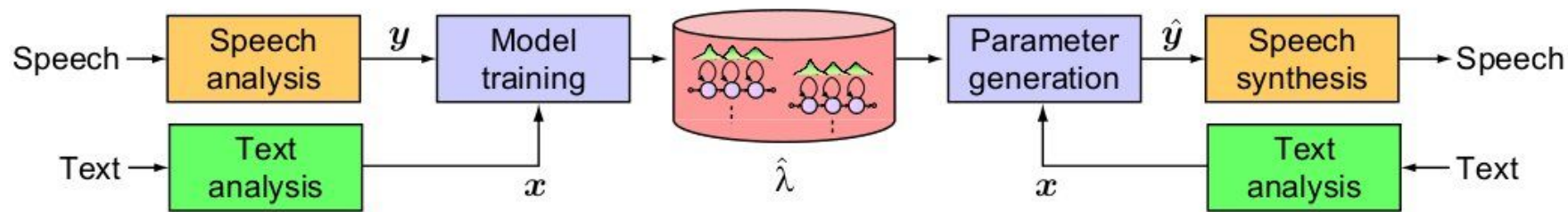
Speech data: 2-10 hours, professional speaker, carefully segmenten and annotated.

Concatenative



— Target cost — Concatenation cost

Statistical Speech Synthesis



Main TTS Technologies

Concatenative speech synthesis + Unit Selection
Concatenate *best* pre-recorded speech *units*

Statistical Parametric Speech Synthesis

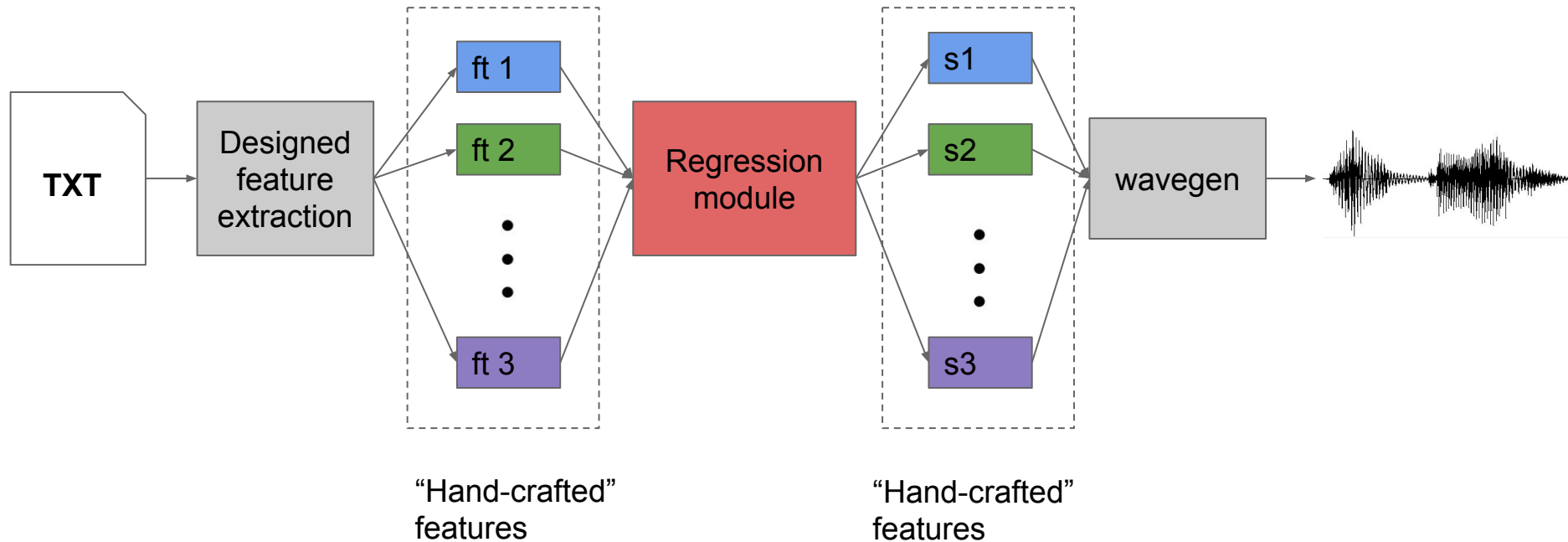
Represent speech waveform using parameters (eg 5ms)
Use statistic generative model
Reconstruct waveform from generated parameters

Hybrid Systems

Concatenative speech synthesis
Select *best* units attending a statistical parametric system

Deep architectures ... but not deep (yet)

Text to Speech: Textual features \rightarrow Spectrum of speech (many coefficients)



Textual features (x)

From text to phoneme (pronunciation)

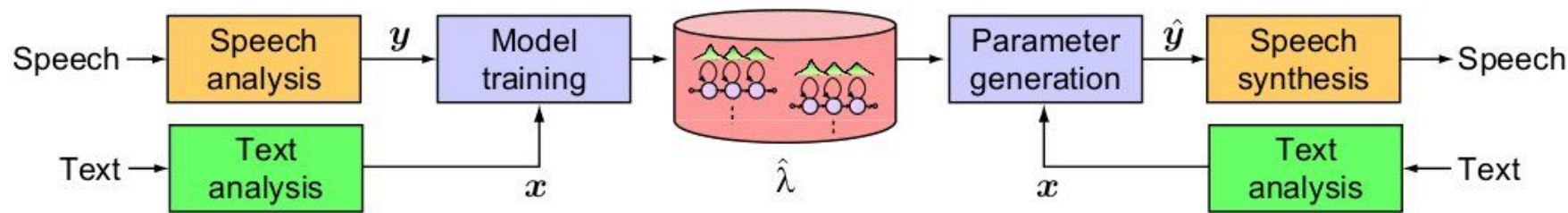
Disambiguation, pronunciation (e.g.: Jan. 26)

From phoneme to phoneme+ (with linguistic features)

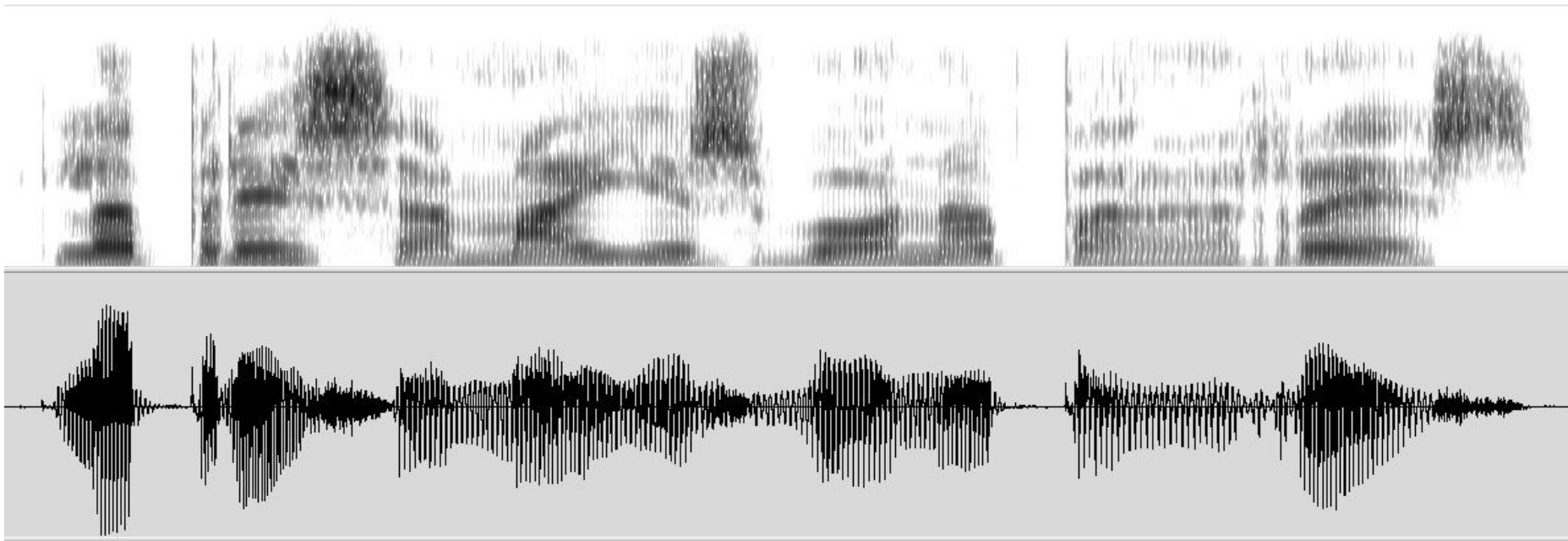
Textual features (x)

- {preceding, succeeding} two phonemes
- Position of current phoneme in current syllable
- # of phonemes at {preceding, current, succeeding} syllable
- {accent, stress} of {preceding, current, succeeding} syllable
- Position of current syllable in current word
- # of {preceding, succeeding} {stressed, accented} syllables in phrase
- # of syllables {from previous, to next} {stressed, accented} syllable
- Guess at part of speech of {preceding, current, succeeding} word
- # of syllables in {preceding, current, succeeding} word
- Position of current word in current phrase
- # of {preceding, succeeding} content words in current phrase
- # of words {from previous, to next} content word
- # of syllables in {preceding, current, succeeding} phrase

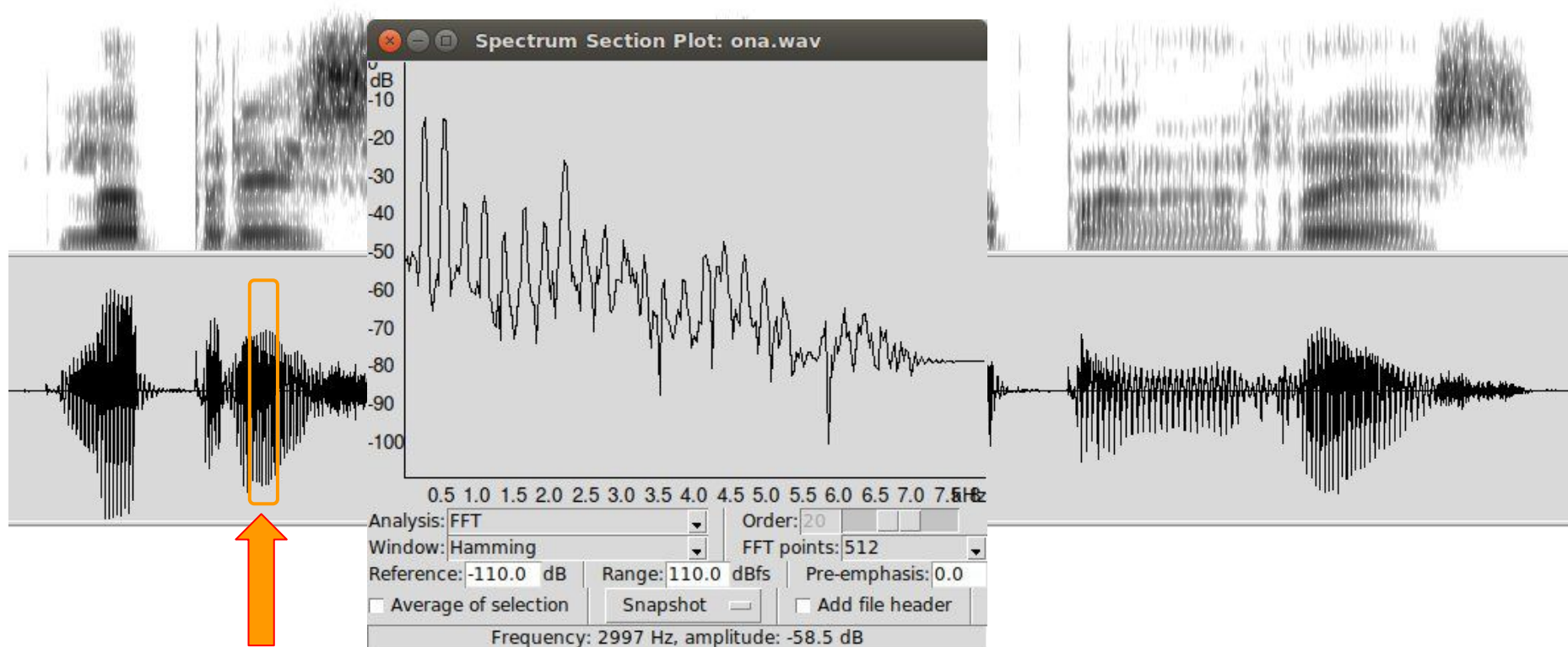
Statistical Speech Synthesis



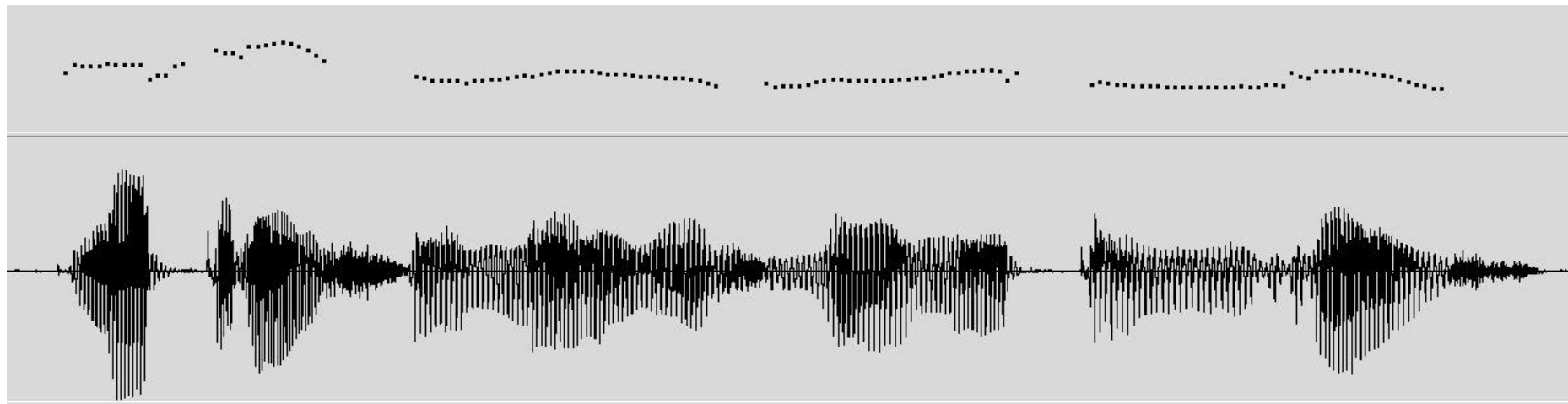
Speech features (y)



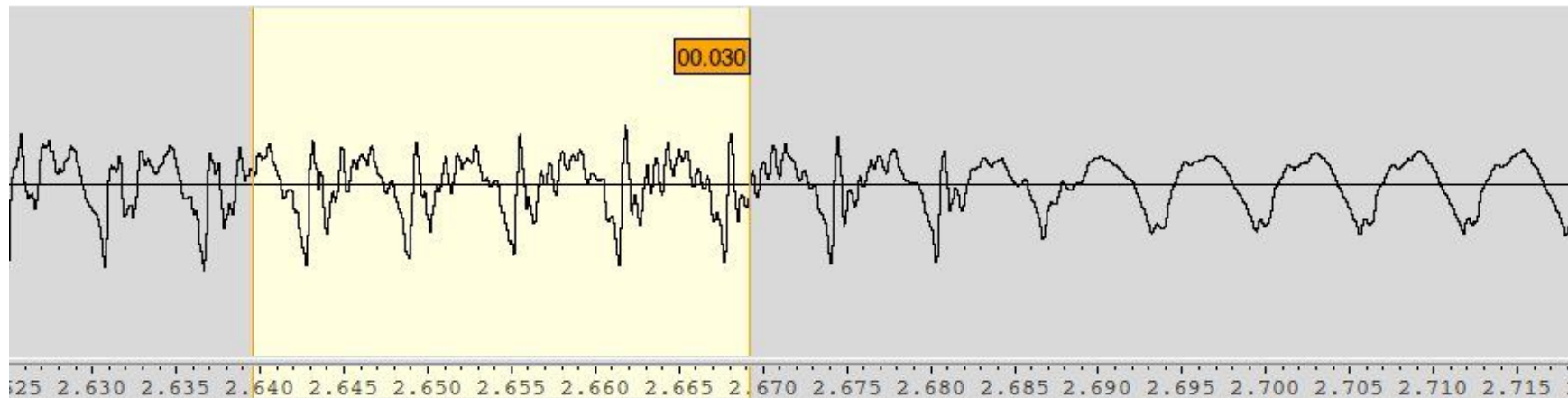
Speech features (y)



Speech features (y)

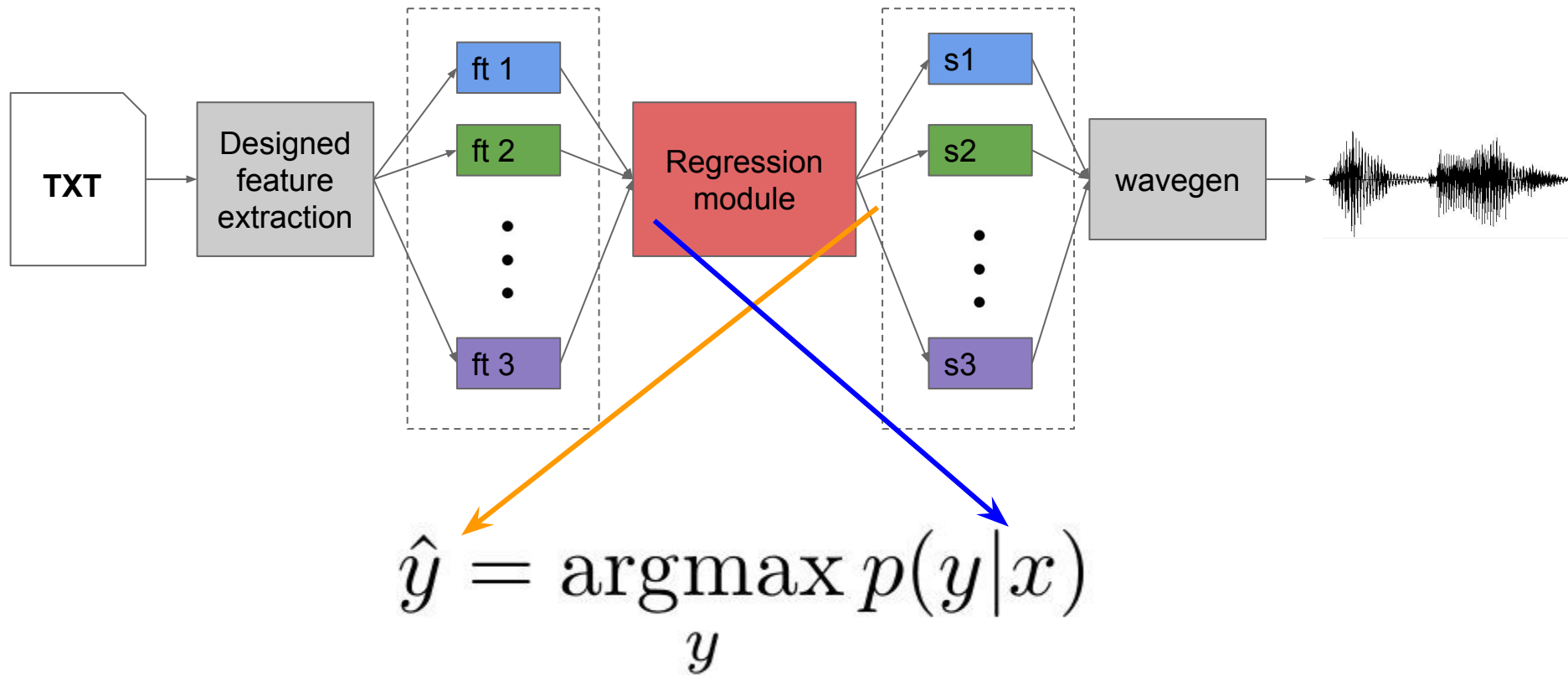


Speech features (y)

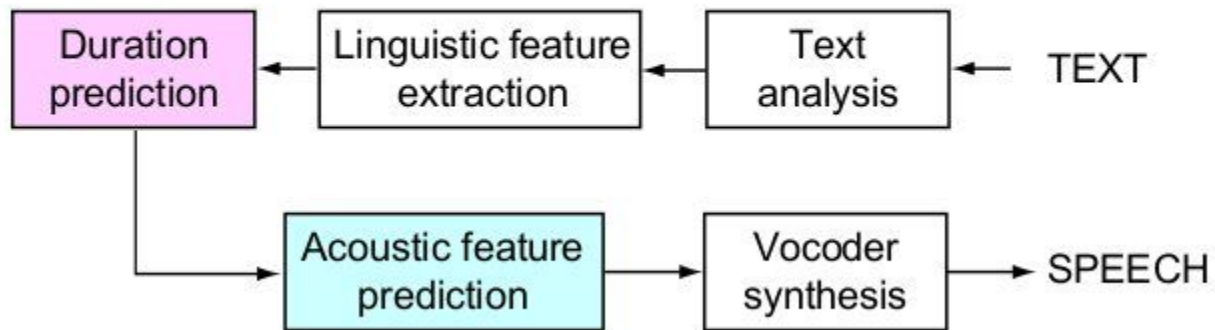


- Rate: ~ 5 ms. (200Hz)
- Spectral features (envelope)
- Excitation features (fundamental frequency, pitch)
- Representation that allows reconstruction: vocoders (Straight, Ahocoder, ...)

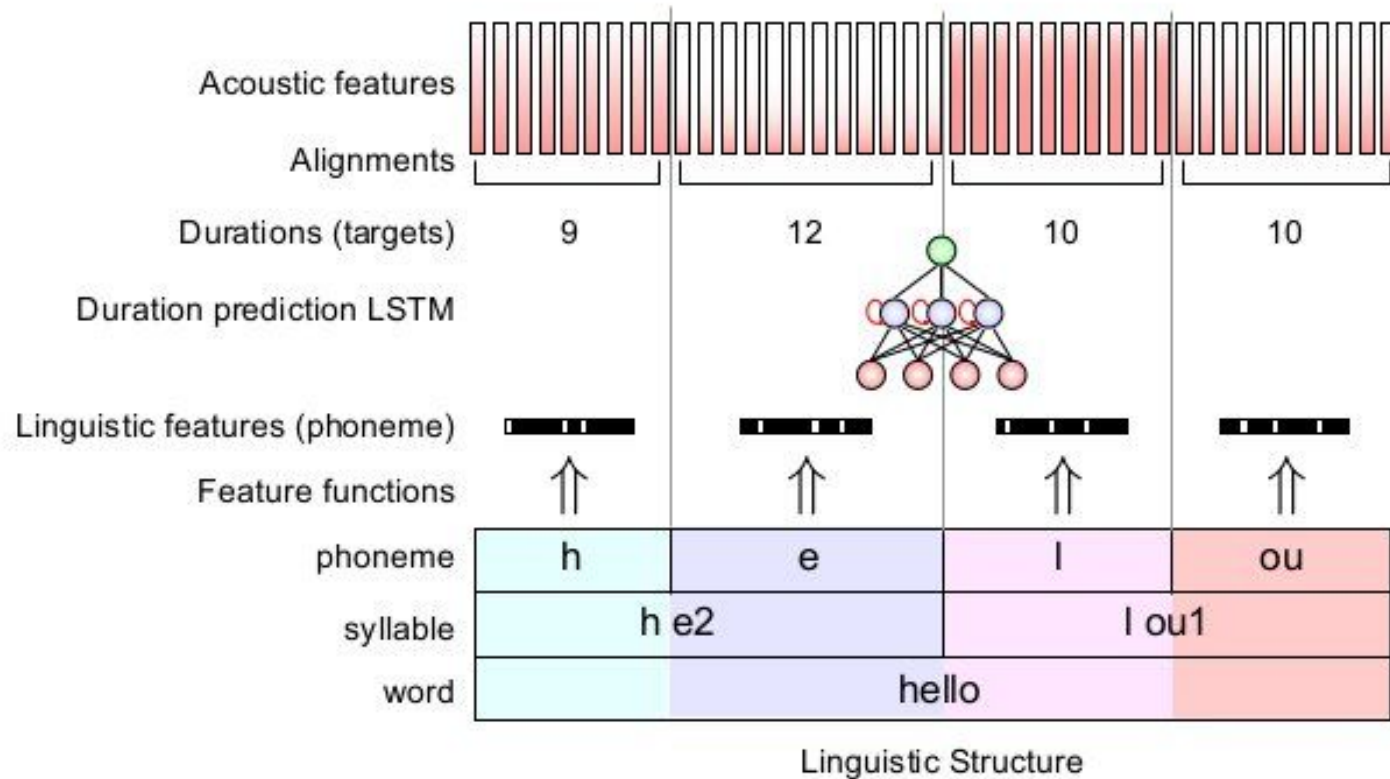
Regression



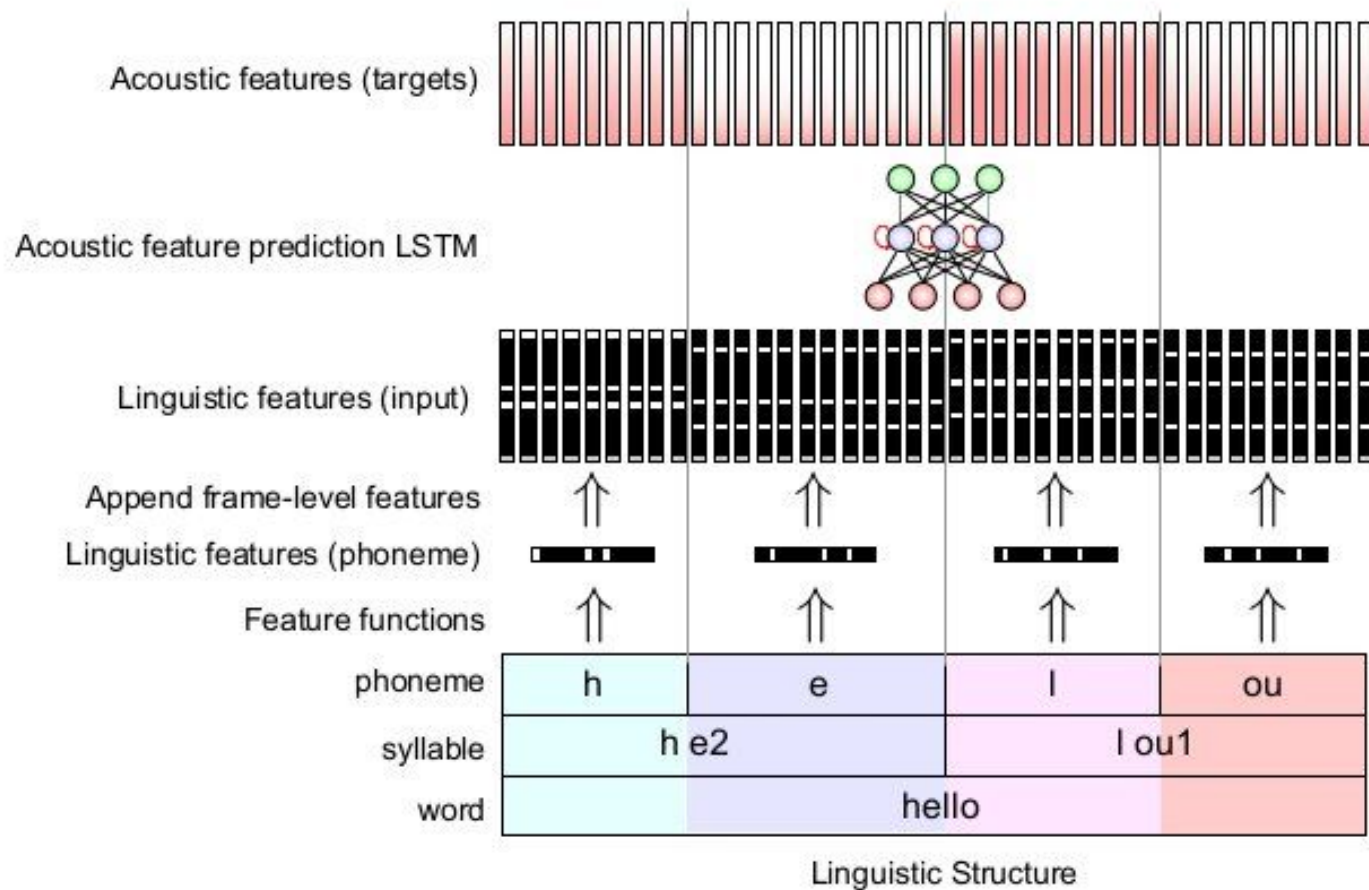
Phoneme rate vs. frame rate



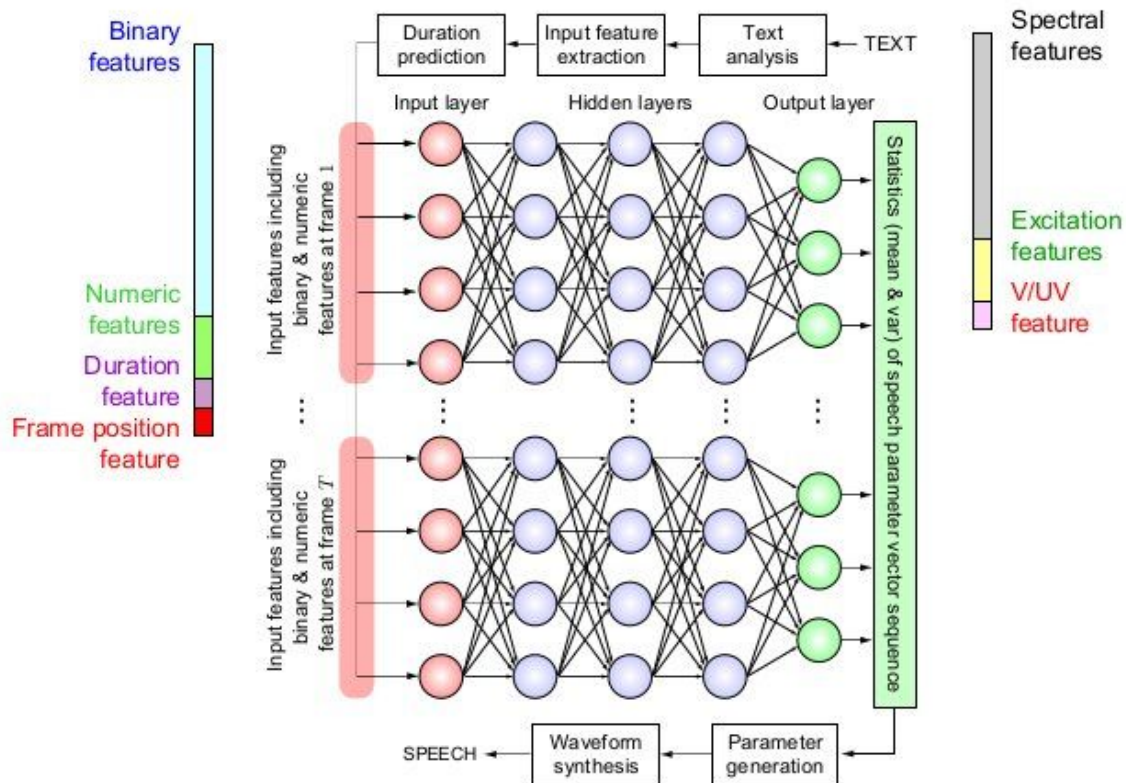
Duration Modeling



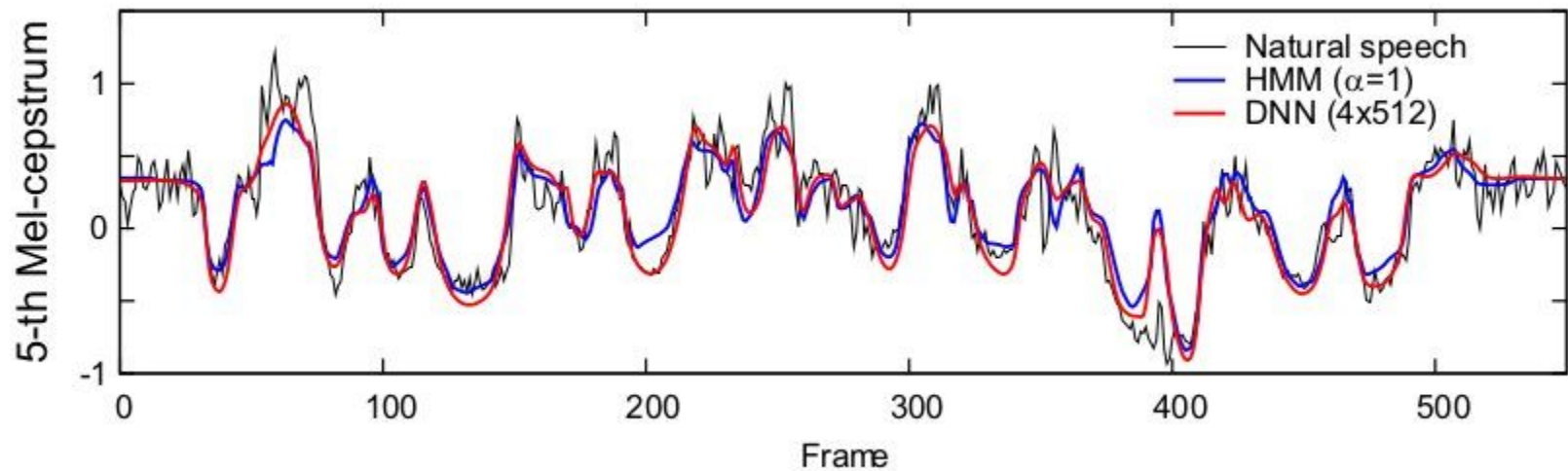
Acoustic Modeling



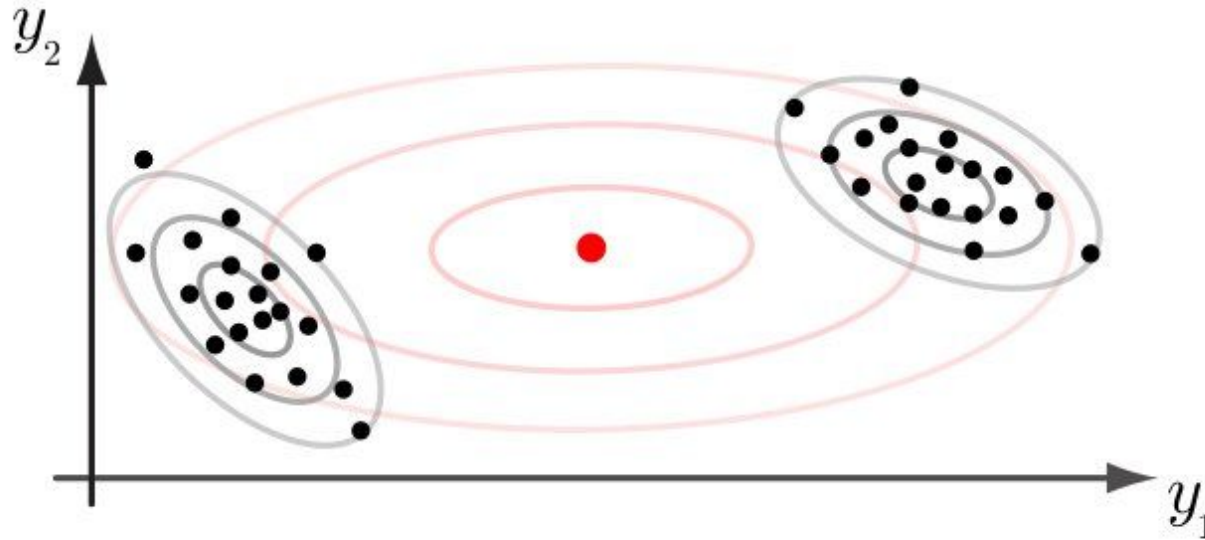
Acoustic Modeling: DNN



Acoustic Modeling: DNN

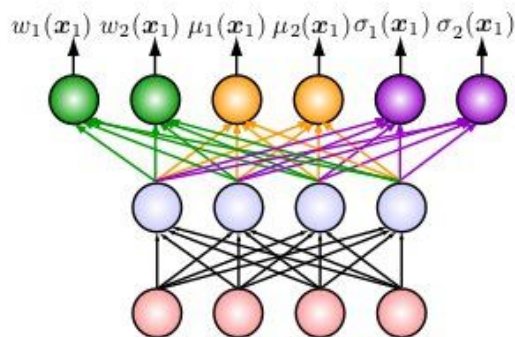
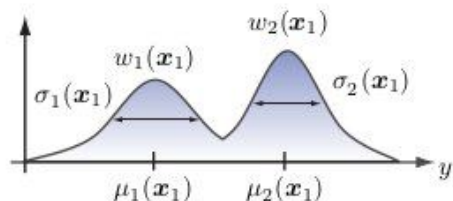


Regression using DNN (problem)



- Data samples
- NN prediction

Mixture density network (MDN)



1-dim, 2-mix MDN

Inputs of activation function

$$z_j = \sum_{i=1}^4 h_i w_{ij}$$

● : Weights → Softmax activation function

$$w_1(\mathbf{x}) = \frac{\exp(z_1)}{\sum_{m=1}^2 \exp(z_m)} \quad w_2(\mathbf{x}) = \frac{\exp(z_2)}{\sum_{m=1}^2 \exp(z_m)}$$

● : Means → Linear activation function

$$\mu_1(\mathbf{x}) = z_3$$

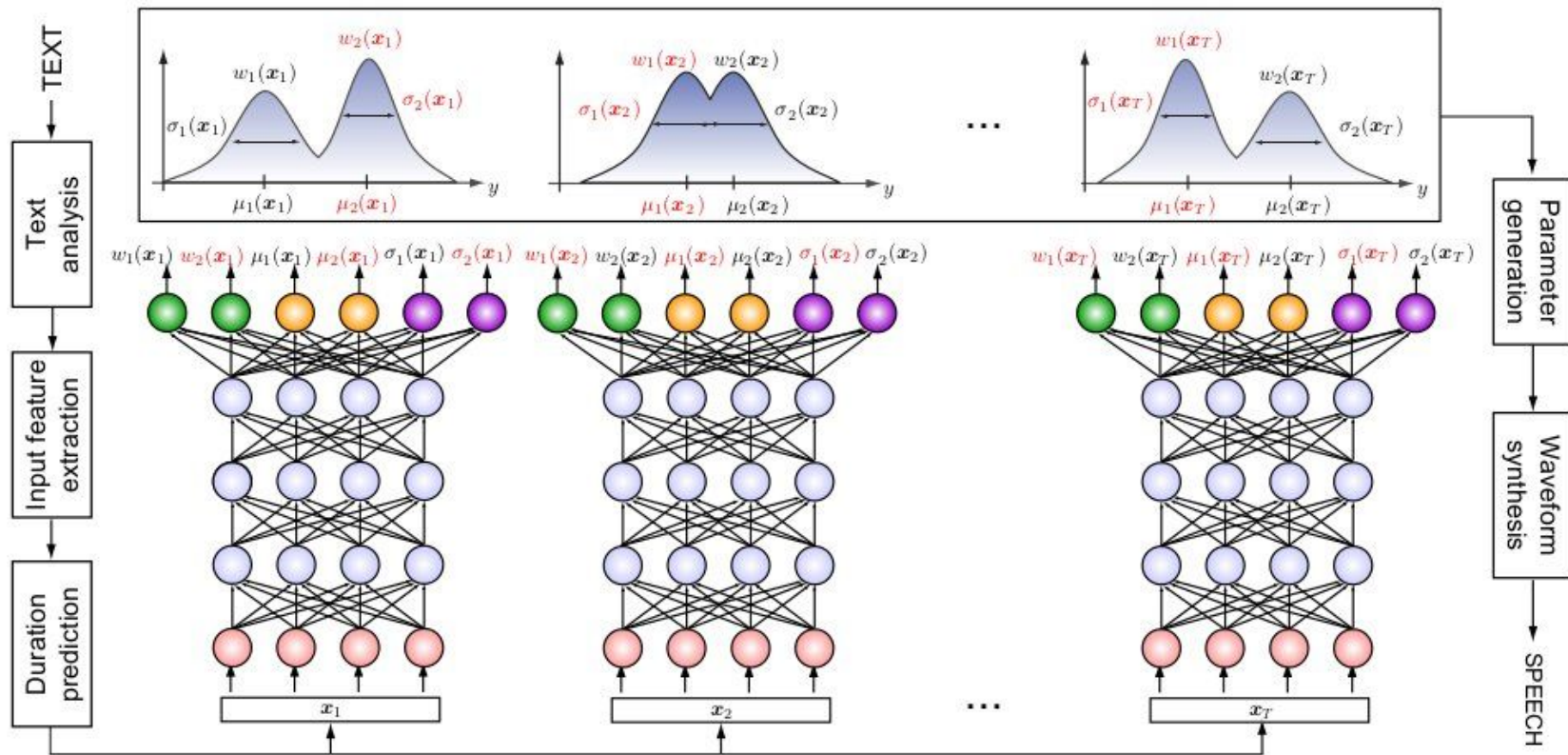
$$\mu_2(\mathbf{x}) = z_4$$

● : Variances → Exponential activation function

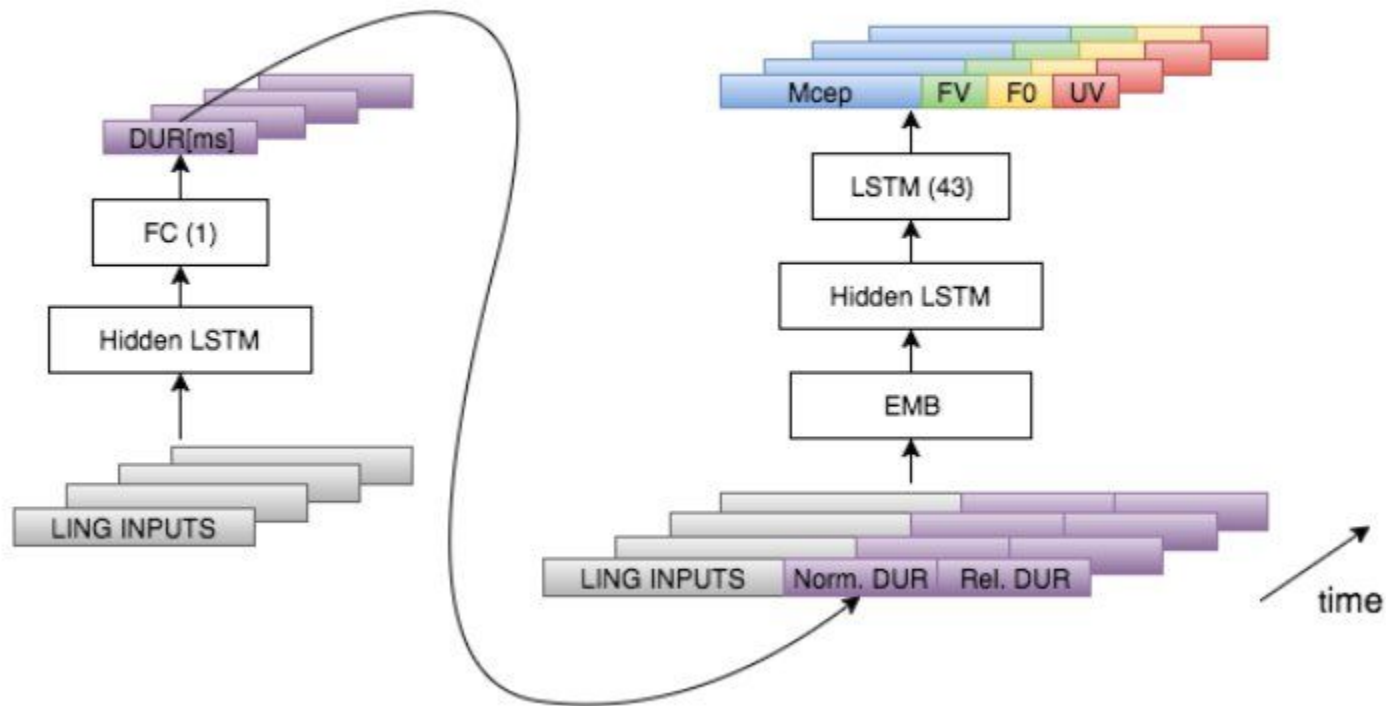
$$\sigma_1(\mathbf{x}) = \exp(z_5)$$

$$\sigma_2(\mathbf{x}) = \exp(z_6)$$

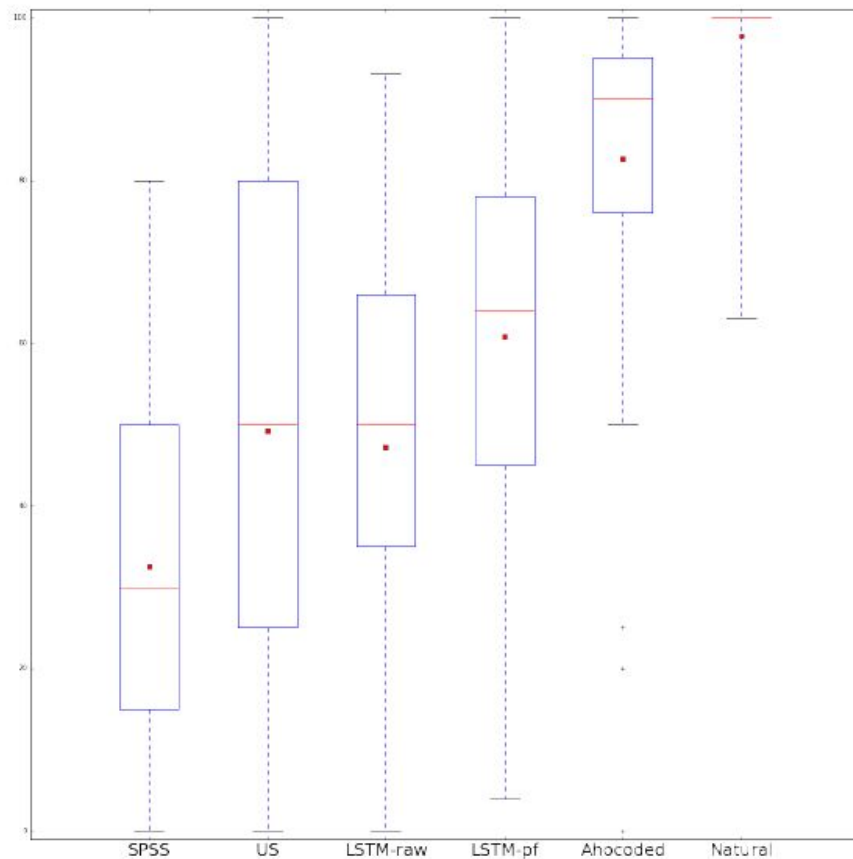
Mixture density network (MDN)



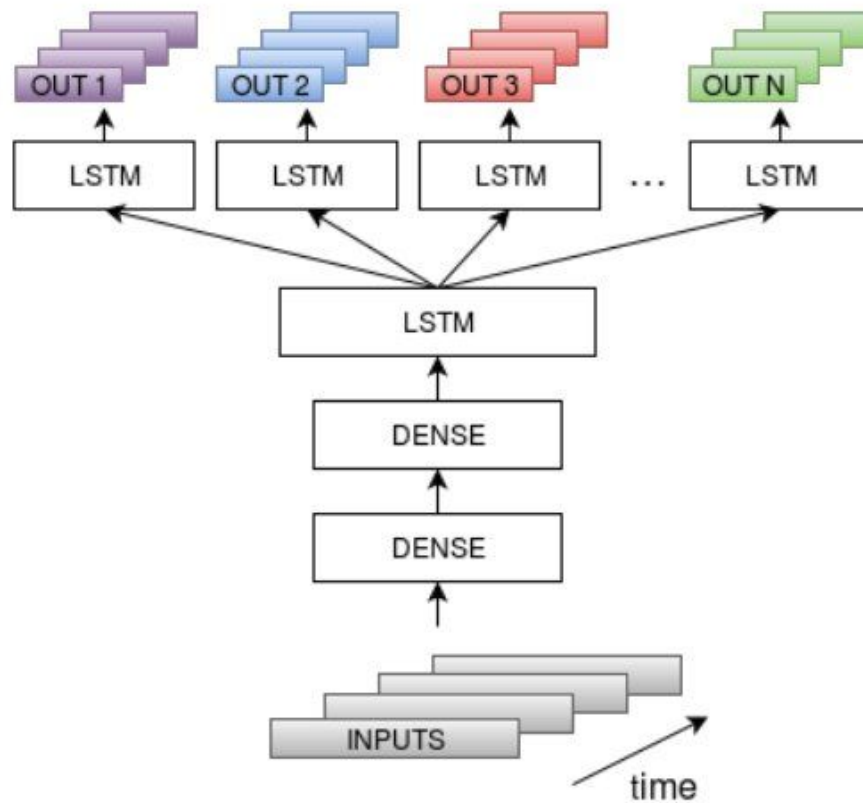
Recurrent Networks: LSTM



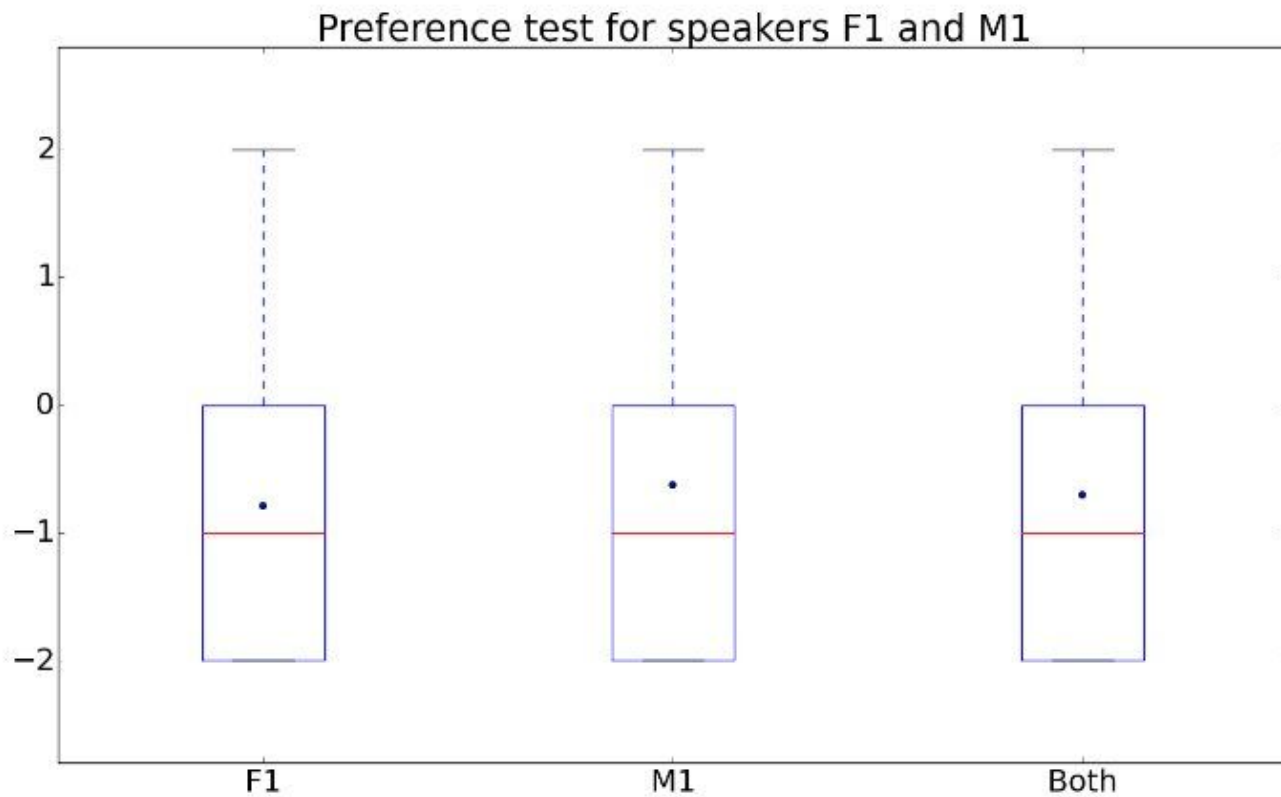
Recurrent Networks: LSTM



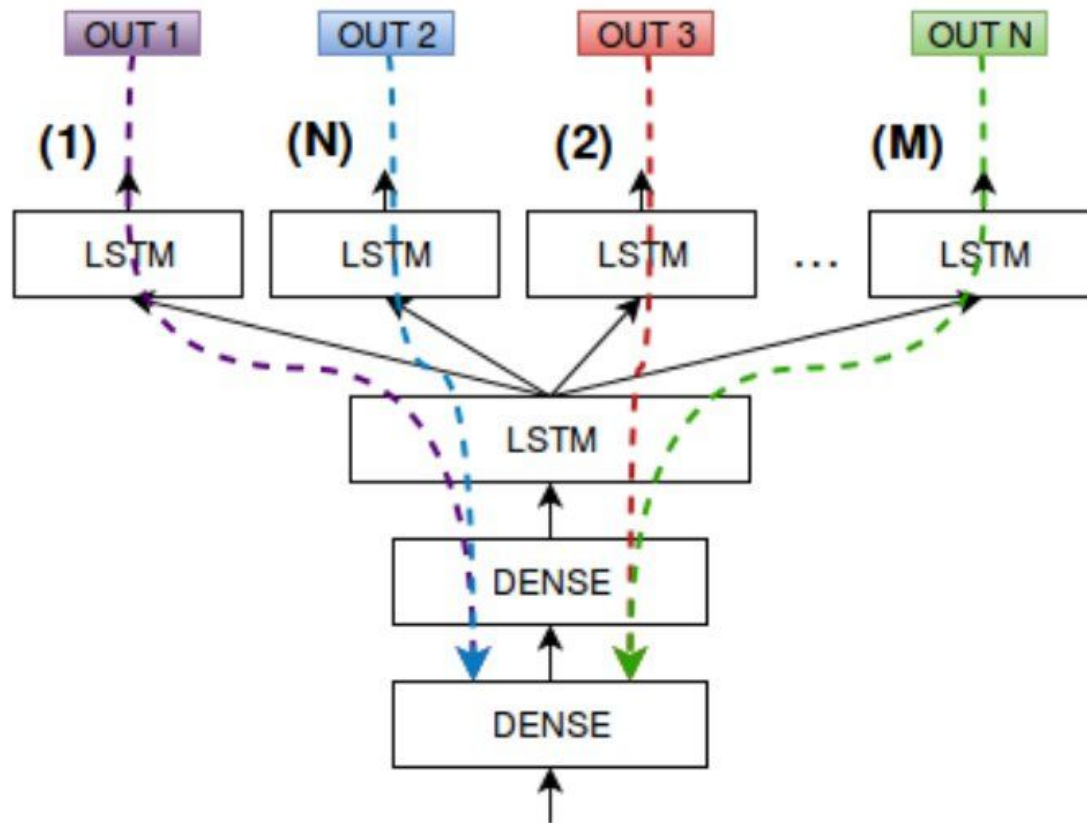
Multi-speaker



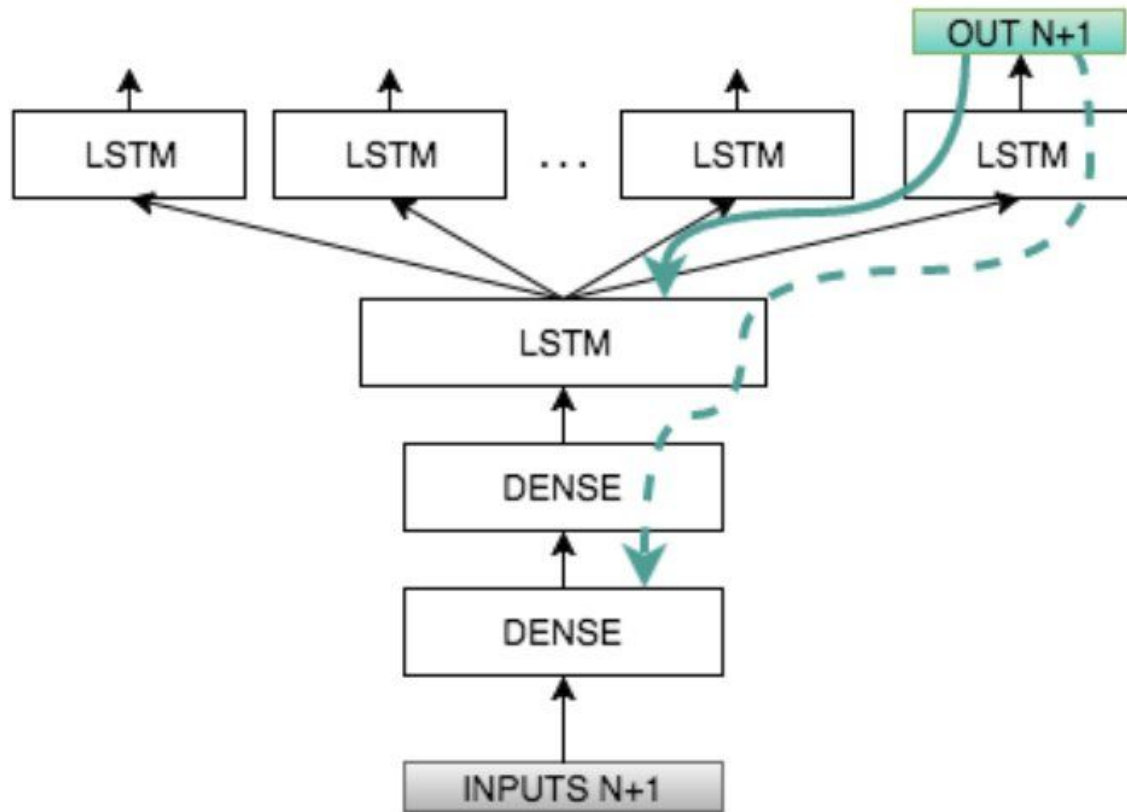
Multi-speaker



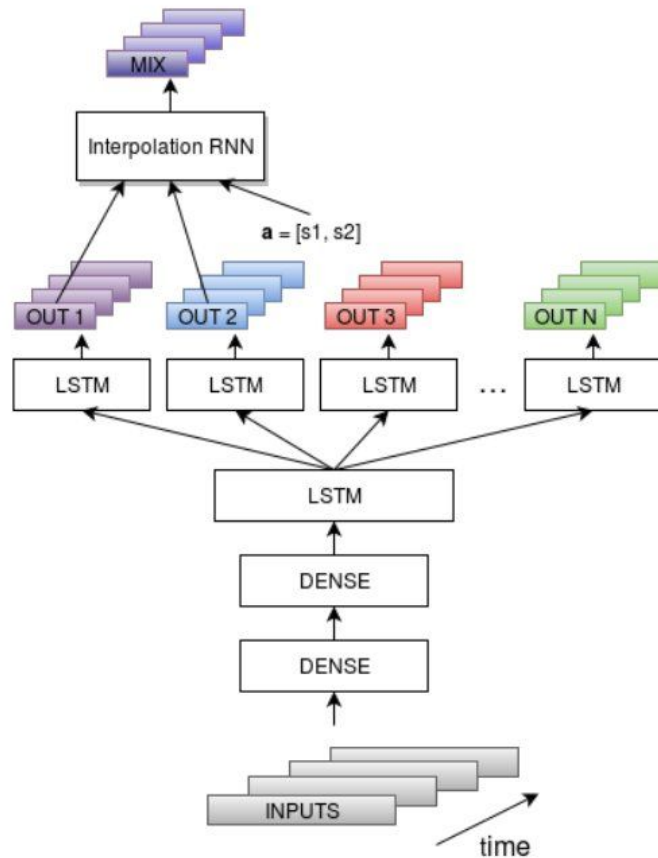
Multi-speaker



Adaptation to new speaker



Speaker interpolation



0 1

0.25 0.75

0.5 0.5

0.75 0.25

1 0

References

Statistical parametric speech synthesis: from HMM to LSTM-RNN.

Heiga Zen, Google

<http://rtthss2015.talp.cat/>

Deep learning applied to Speech Synthesis, Msc Thesis

Santiago Pascual, UPC