# Exploring the multidimensional representation of unidimensional speech acoustic parameters extracted by deep unsupervised models

M. Jacquelin[1,2], M. Garnier[1], L. Girin[1], R. Vincent[2], O. Perrotin[1]

[1] GIPSA-Lab, Université Grenoble Alpes    [2] Vogo
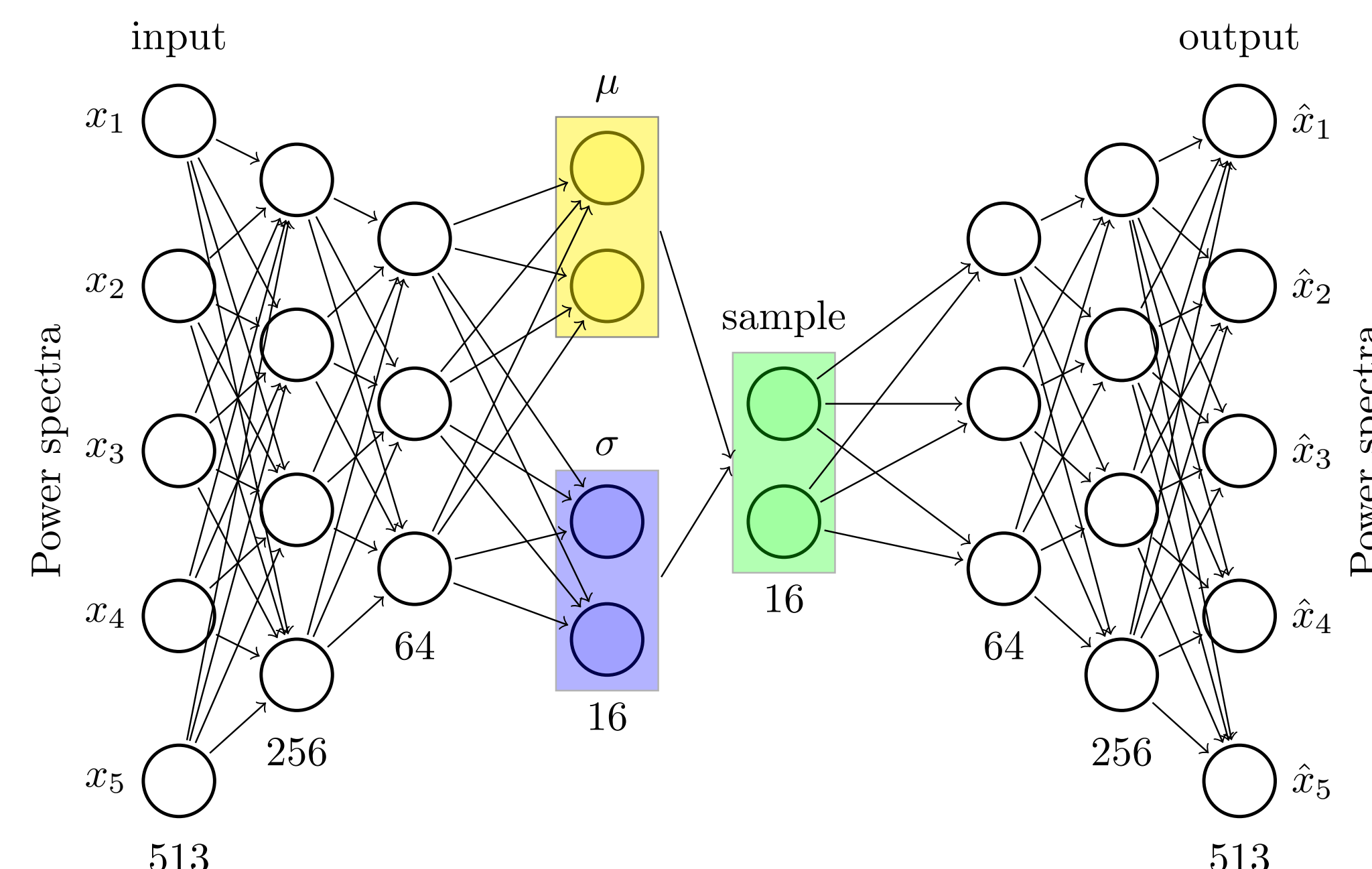
## Context & problematic

The work of [1] demonstrate that the fundamental frequency ($f_0$) and the frequency of the first three formants ($F_{1,2,3}$) are encoded in multiples dimensions in the latent space of unsupervised models. This raises the following questions:

- Why unsupervised models need multiple dimensions to encode acoustic parameters ?
- What type of information or acoustic variability is captured by each of these latent dimensions ?
- Can we control the latent space of our model to transform the variability of the acoustic parameters ?

## Model used

The variational autoencoder architecture used in this work is similar to that used in [1]:

- The model was trained on 20 hours of VCTK on speakers not used for testing;
- The loss function is the weighted sum of the Itakura-Saito divergence;



## Datas & analysis methods

**Datasets:**

- A natural speech test set called $D_{\text{NS},x}^{\text{test}}$ was created by extracting 3 hours from VCTK [2] ;
- We have built a synthetic speech test set called $D_{\text{SS},x}^{\text{test}}$ with Soundgen [3];
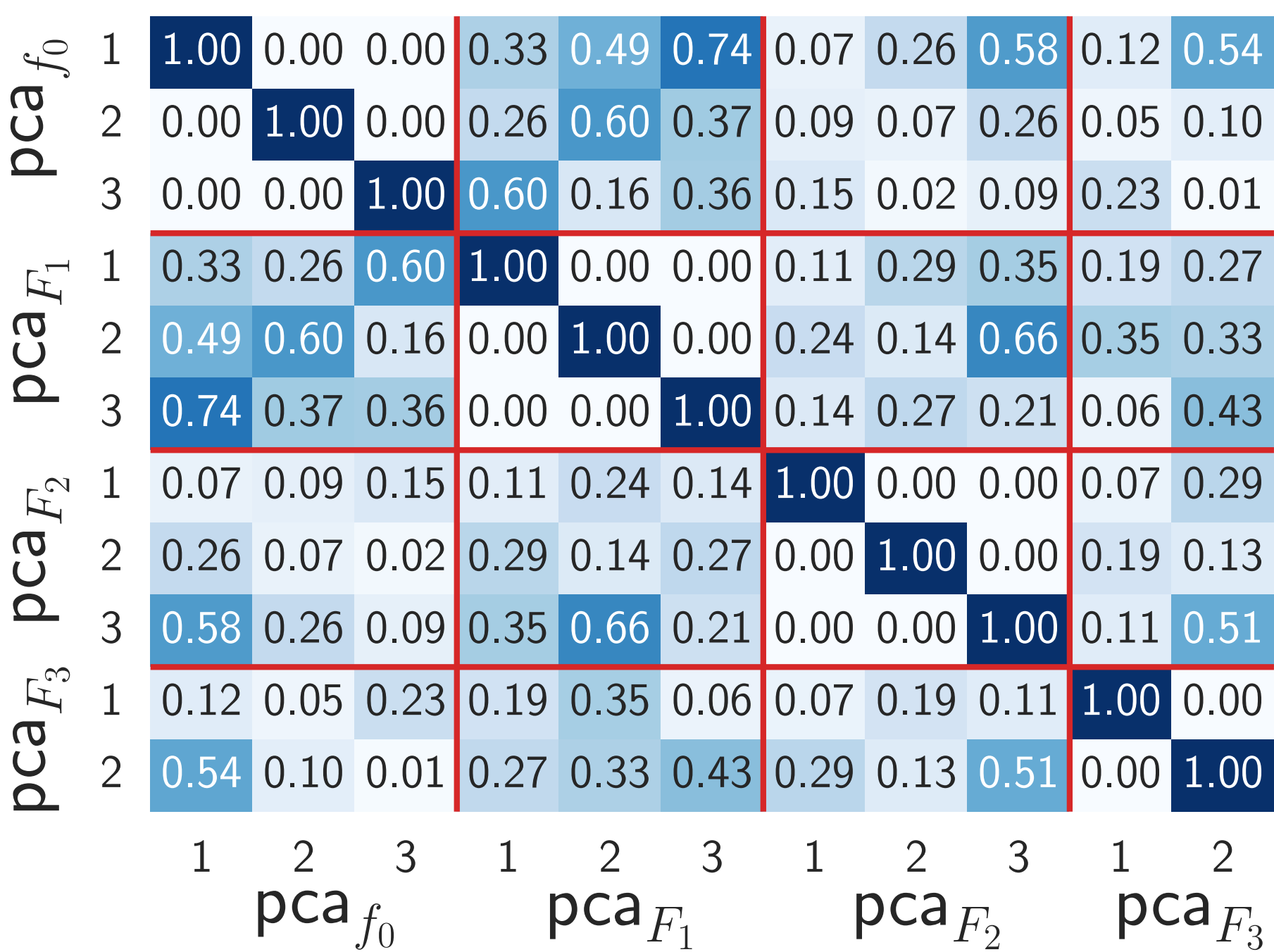- Acoustic parameters $\mathcal{F} \in \{f_0, F_1, F_2, F_3\}$;

**Methods:**

- Principal Component Analysis (PCA) maximizes the variance of the projected data;
- Linear Regression (LR) estimate the relationship between variables that can be separated by classes ($\mathcal{C}$);

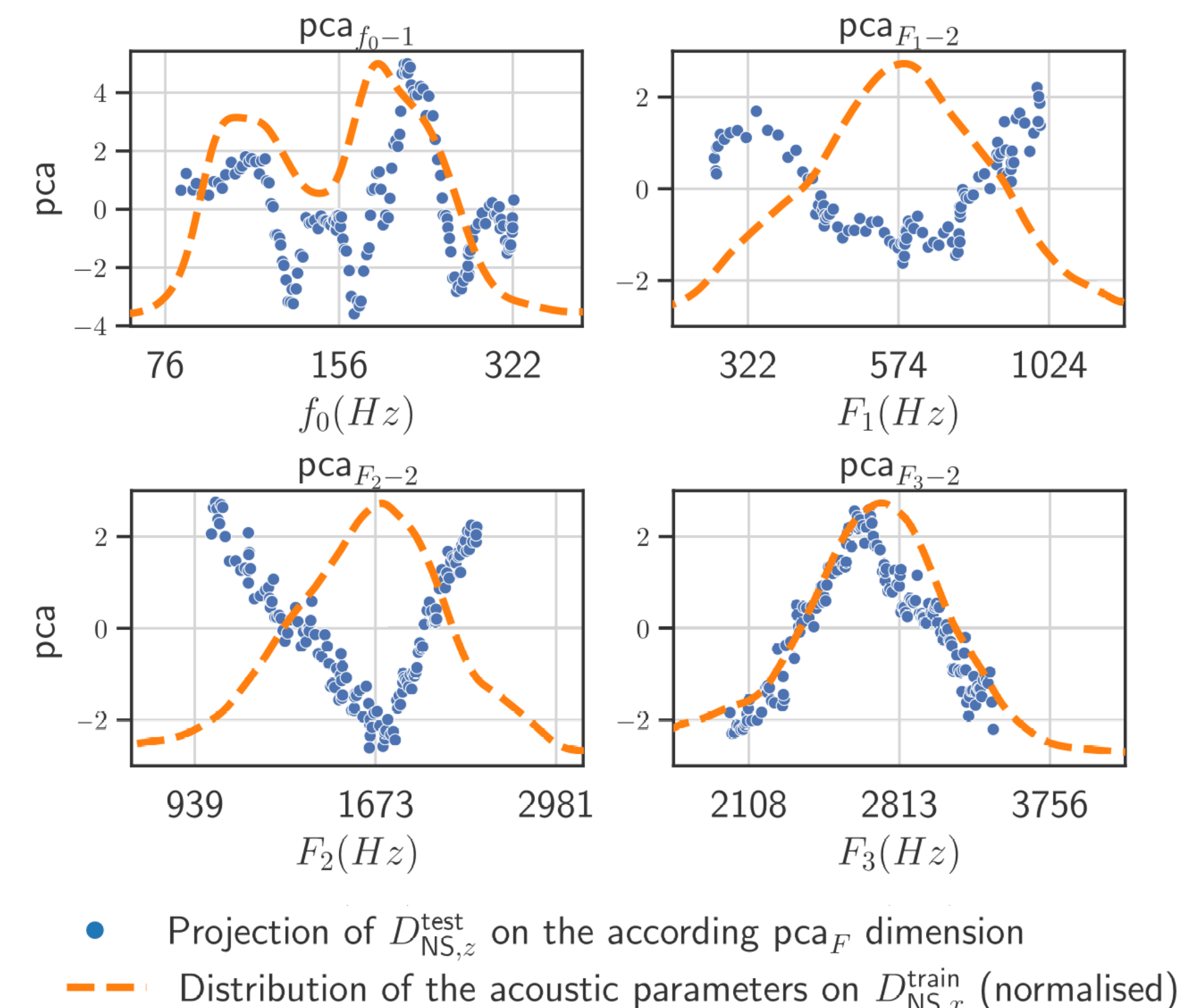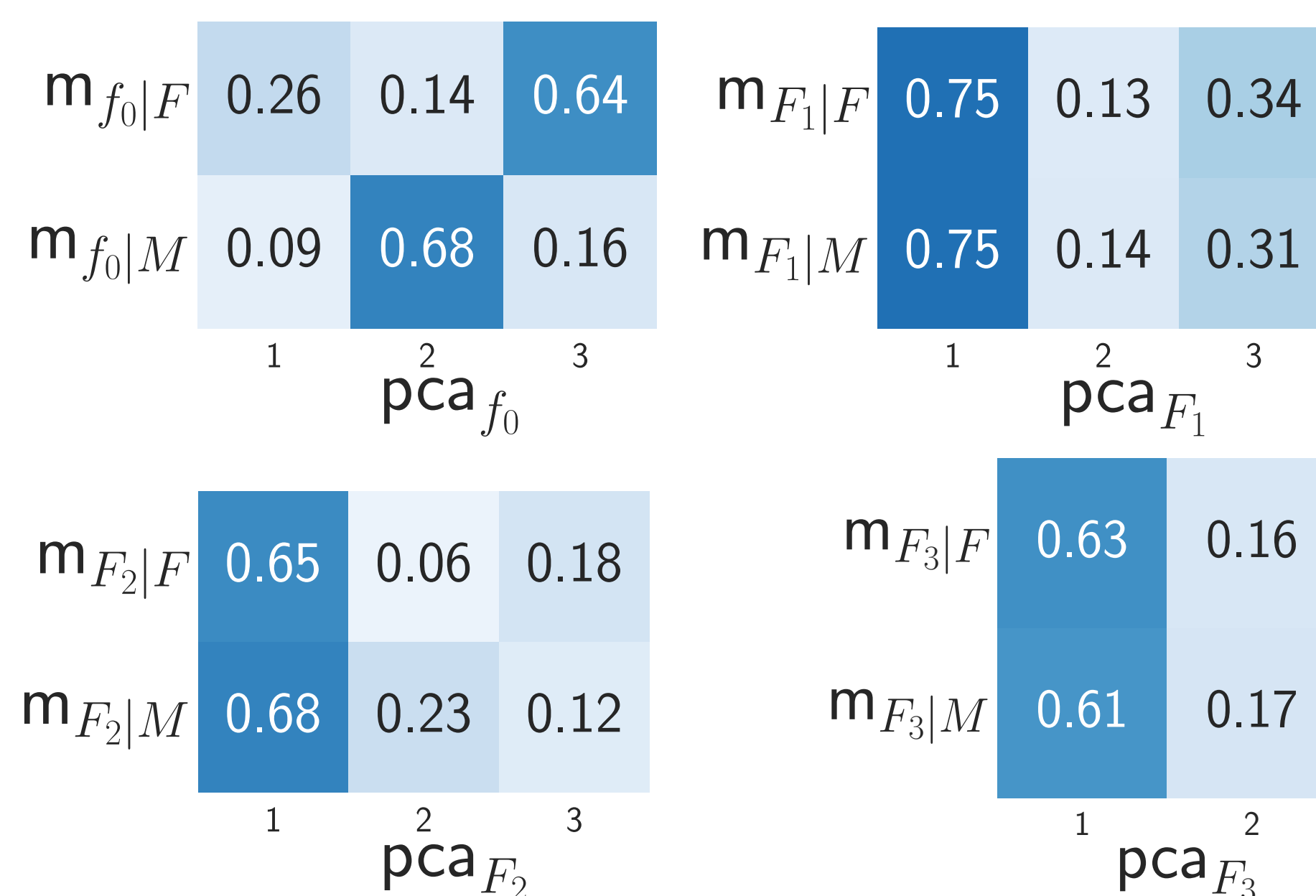| Signal | Encoding | Latent | Analysis | Directions |
|---|---|---|---|---|
| $D_{\text{SS}(\mathcal{F}),x}^{\text{test}}$ | VAE | $D_{\text{SS}(\mathcal{F}),z}^{\text{test}}$ | PCA | $\text{pca}_\mathcal{F}$ |
| $D_{\text{NS},x}^{\text{test}}$ | VAE | $D_{\text{NS},z}^{\text{test}}$ | LR($\mathcal{F}$) | $\text{m}_\mathcal{F}$ |
| | | | LR($\mathcal{F}|\mathcal{C}$) | $\text{m}_{\mathcal{F}|\mathcal{C}}$ |

## Results

- **The variation of each acoustic parameter is encoded by multiple dimensions in the latent space.**
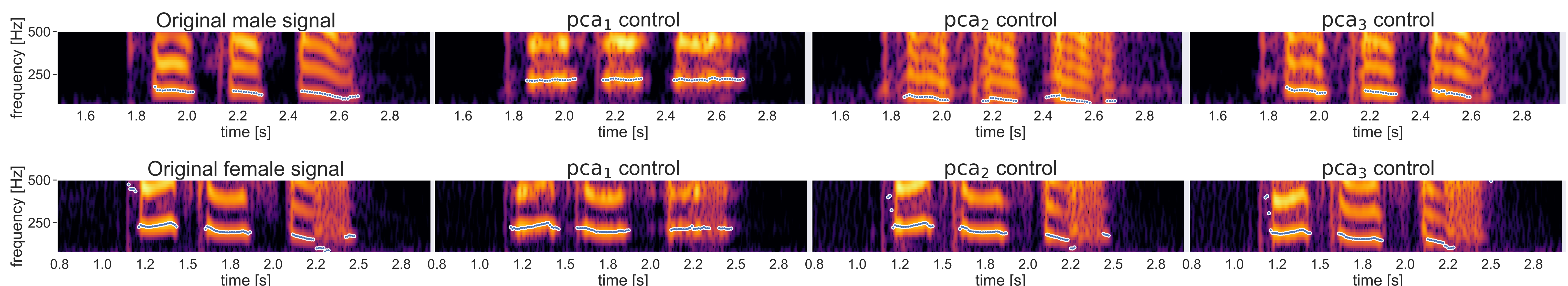


*Can we relate these dimensions to variations of acoustic parameters?*

- **The multidimensional representation of acoustic parameters is closely related to the multimodality of the parameter distribution.**





- Projection of $D_{\text{NS},z}^{\text{test}}$ on the according $\text{pca}_F$ dimension
- Distribution of the acoustic parameters on $D_{\text{NS},x}^{\text{train}}$ (normalised)

- **The variation of acoustic parameter can be controlled directly from the encoded latent space.**



## Conclusion

- ✔ **Identification of the directions that best explain the variation of selected acoustic features**
- ✔ **Highlighting the link between multimodality of parameter distribution and multidimensional representation**
- ✔ **Demonstrate the ability to control the variation of acoustic parameter with the encoded latent space**

## References

[1] S. Sadok, S. Leglaive, L. Girin, X. Alameda-Pineda, R. Séguier, *Learning and controlling the source-filter representation of speech with a variational autoencoder*, in Speech Communication, 2023, vol. 148, pp. 53-65.

[2] J. Yamagishi, C. Veaux, K. MacDonald, CSTR VCTK corpus: English multi-speaker corpus for CSTR voice cloning toolkit, 2019.

[3] A. Anikin, *Soundgen: an open-source tool for synthesizing non-verbal vocalizations*, in Behavior research methods, 2019, vol. 51, pp. 778–792.

**Grenoble Images Parole Signal Automatique**

UMR CNRS 5216 - Grenoble Campus

38400 Saint Martin d'Hères - FRANCE