

Hybrid deep learning for audio

Gaël RICHARD*

Professor, Telecom Paris, Institut polytechnique de Paris

*work with collaborators and in particular **K. Schulze-Forster**, C. Doire, R. Badeau

Content

- Context and motivation
- Towards hybrid deep learning
 - Some examples in other domains
 - Hybrid deep learning in audio
 - A specific example in unsupervised source separation
- Discussion and conclusion

Context and motivation

- Machine learning: a growing trend towards pure “Data-driven” deep learning approaches
- High performances but some main limitations:
 - *“Knowledge” is learned (only) from data*
 - *Complexity: overparametrized models (> 100 millions parameters)*
 - Overconsumption regime
 - Non-interpretable/non-controllable

Context and motivation

- Machine learning: a growing trend towards pure “Data-driven” deep learning approaches
- High performances but some main limitations:
 - “*Knowledge*” is learned (only) from data
 - *Complexity: overparametrized models (> 100 millions parameters)*
 - Overconsumption regime
 - Non-interpretable/non-controllable
- **The main goal of my ERC project Hi-Audio :**



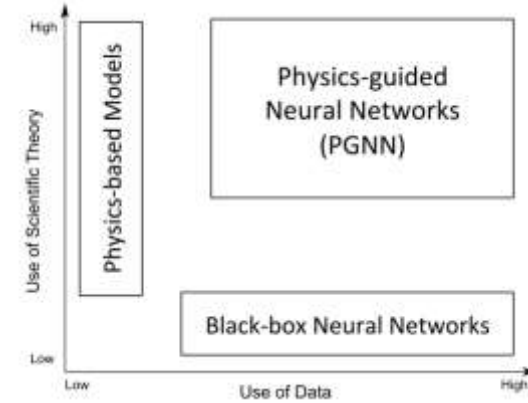
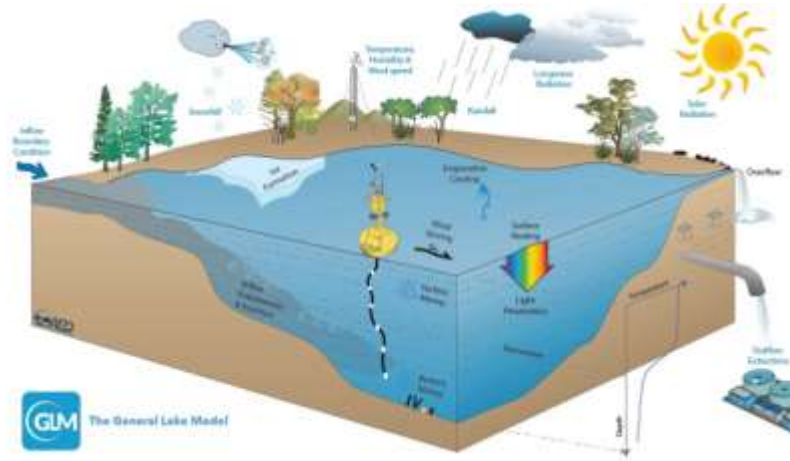
Main goal : To build controllable and frugal machine listening models based on expressive generative modelling

My approach: to build *Hybrid deep learning models*, by **integrating our prior knowledge** about the nature of the processed data.

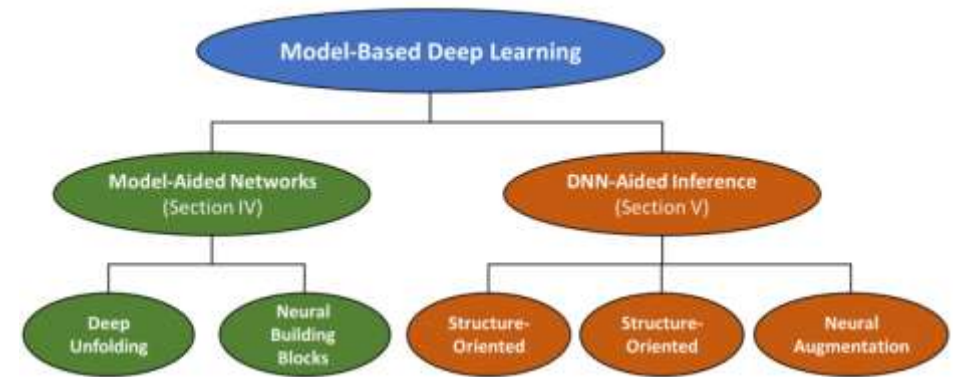
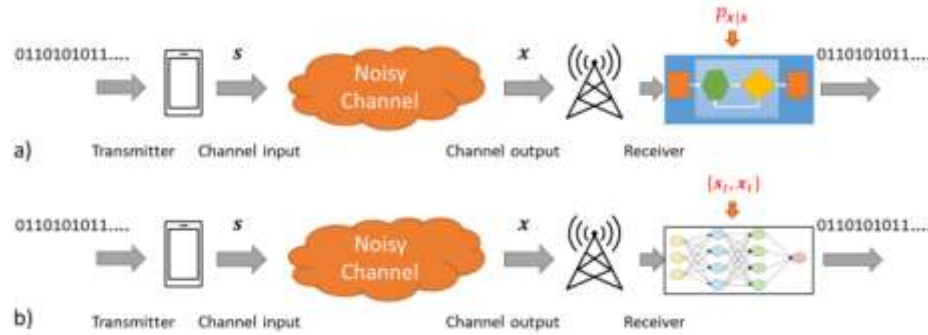
Towards Hybrid deep learning

... some prior works.

- Physics-guided neural networks in remote sensing [1],



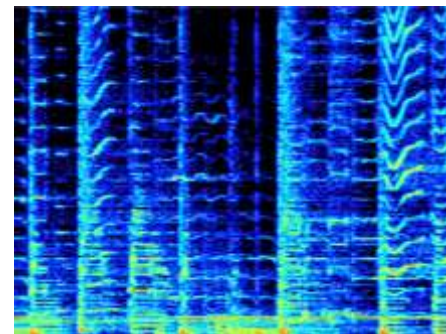
- Digital communication and Image restoration [2,3]



[1] A. Karpatne & al. "Physics-guided Neural Networks (PGNN): An Application in Lake Temperature Modeling," arXiv, 1710.11431, 2017.
 [2] B. Lecouat & al., "Fully Trainable and Interpretable Non-Local Sparse Models for Image Restoration.," 2020. (hal-02414291v2).
 [3] N. Shlezinger, & al., "Model-Based Deep Learning," arXiv, 2012.08405, 2020.

Why Hybrid deep learning is interesting for audio ?

- We often have some good apriori knowledge of the sources
- We have a long history of audio signal models:
 - **Audio perception:** hearing model, psychoacoustics,...
 - **Audio sound production:** source-filter, periodic/aperiodic, physical model,...
 - **Audio propagation:** room acoustics, reverberation, ...
 - **Audio signal models:** sparsity, factorisation, time-frequency representations, decomposition models ...
- Audio requires specific deep neural networks architectures (compared to image processing)

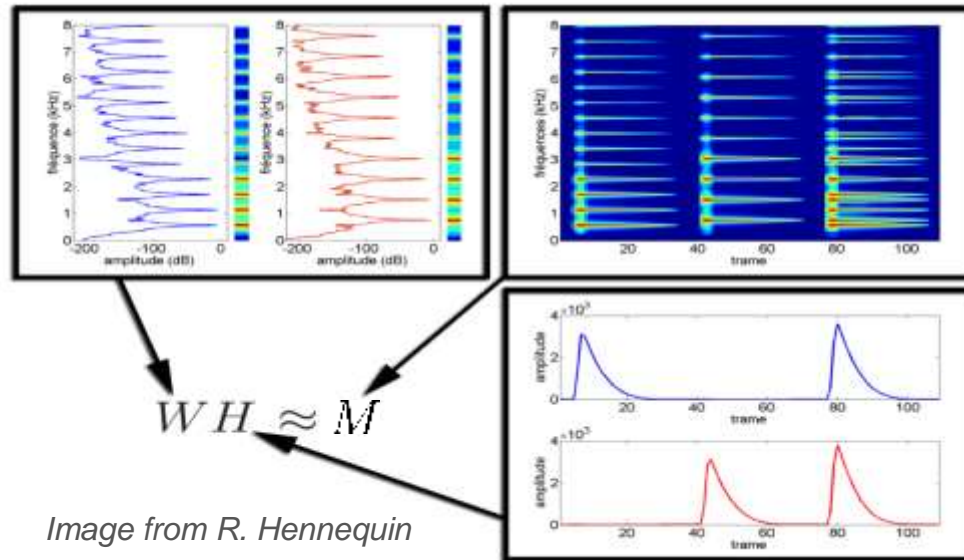


Towards Hybrid deep learning

... some prior works in Audio.

- **Signal models can be used as an advanced representation:**
 - An example: non-negative factorization models with CNNs for audio scene classification

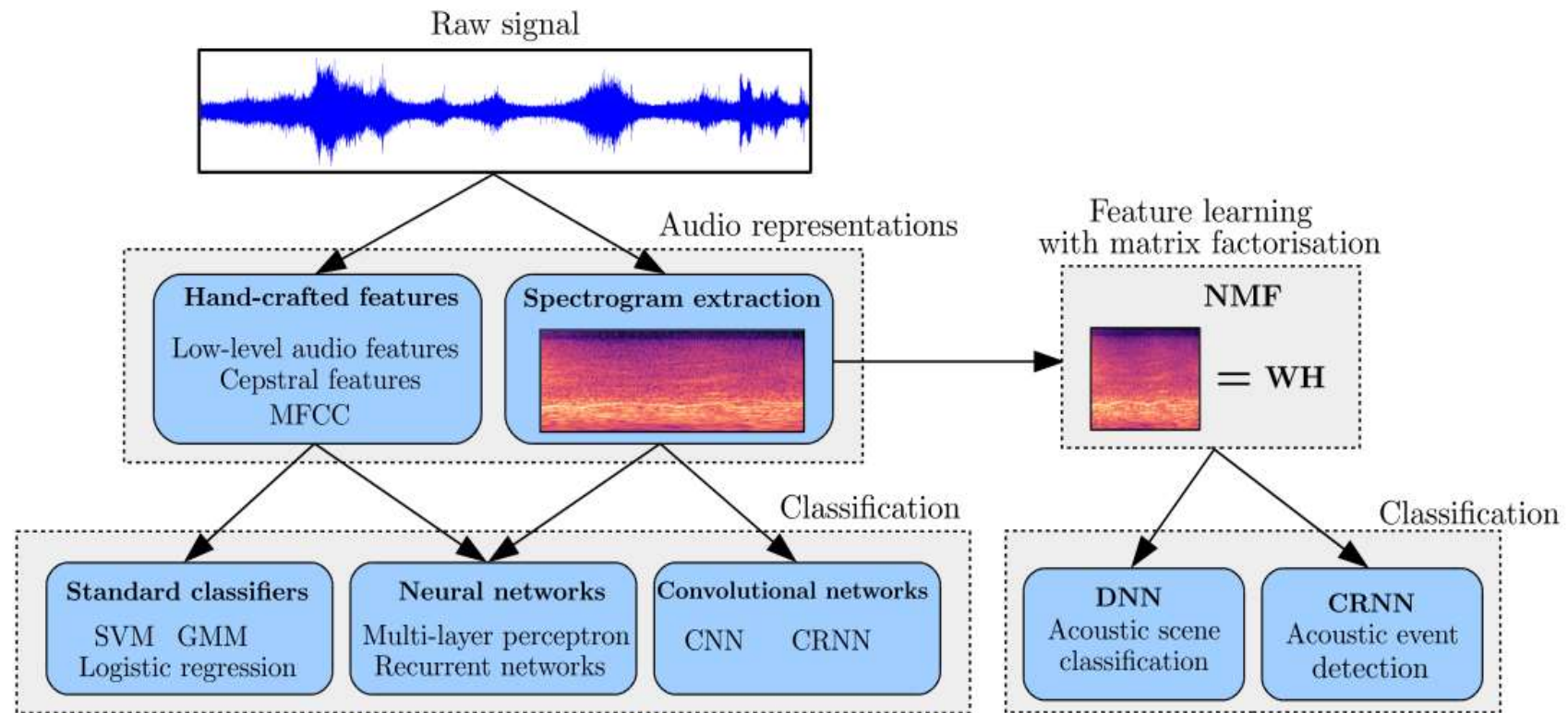
Principle of Non-Negative Matrix Factorization on Audio spectrograms



Towards Hybrid deep learning

... some prior works in Audio.

- Feature learning with NMF for audio scene classification



V. Bisot & al., "Feature Learning with Matrix Factorization Applied to Acoustic Scene Classification", ACM/IEEE Trans. on ASLP, vol. 25, no. 6, 2017

V. Bisot & al., Leveraging deep neural networks with nonnegative representations for improved environmental sound classification IEEE International Workshop on Machine Learning for Signal Processing MLSP, Sep 2017, Tokyo,

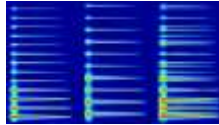


Towards Hybrid deep learning

... some prior works in Audio.

- Deep NMF : the concept of deep unrolling

- Classic NMF



$$\mathbf{M} \approx \mathbf{W}\mathbf{H} = \hat{\mathbf{M}}$$

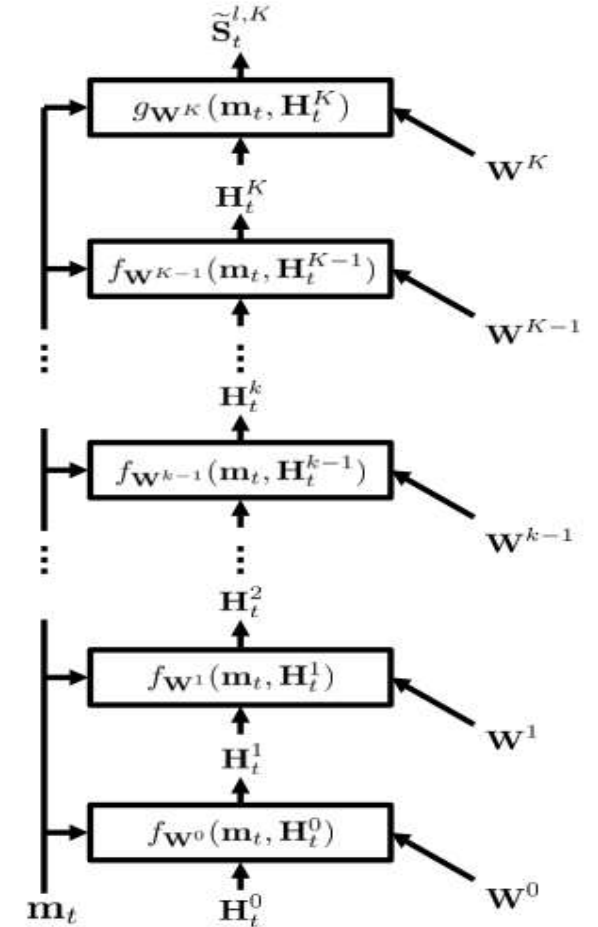
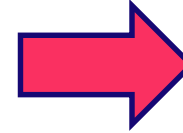
- Minimizing a distance

$$D(\mathbf{M}, \hat{\mathbf{M}}) = \sum_{f=1}^F \sum_{n=1}^N d(v_{fn} | \hat{v}_{fn})$$

- ..towards iterative update rules

$$\mathbf{H} \leftarrow \mathbf{H} \otimes \frac{\mathbf{W}^T \mathbf{M}}{\mathbf{W}^T (\mathbf{W}\mathbf{H})}$$

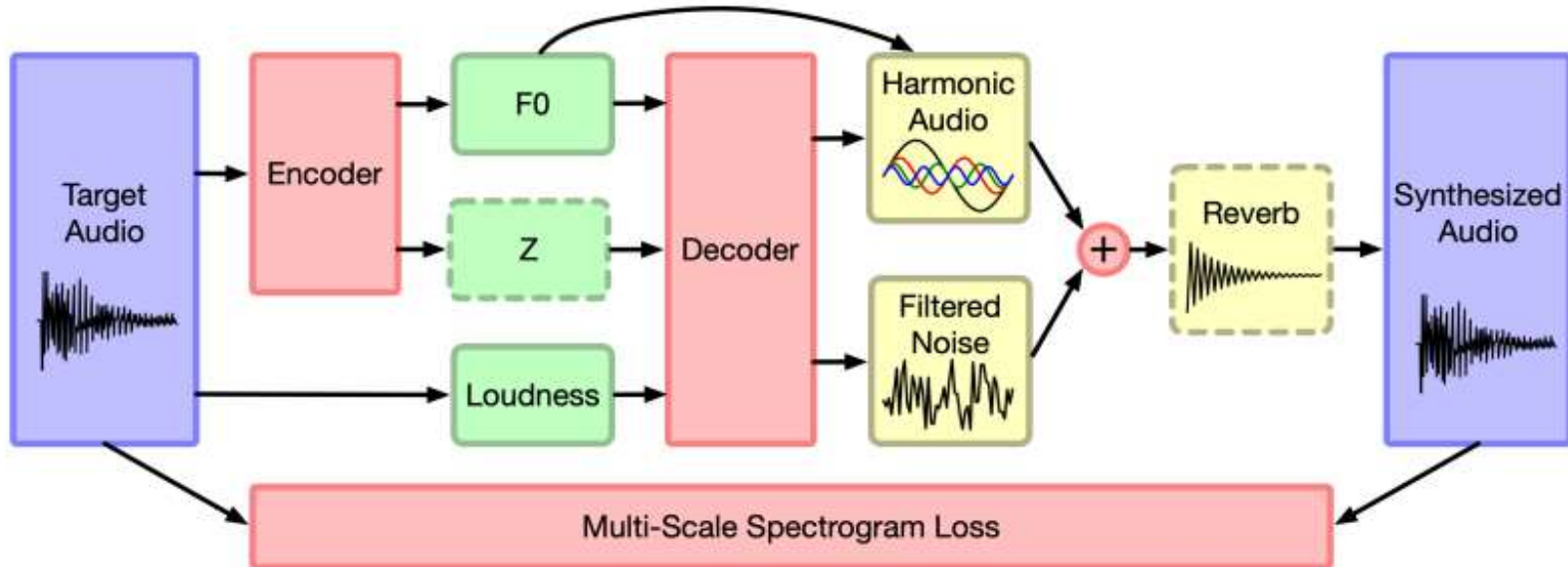
$$\mathbf{W} \leftarrow \mathbf{W} \otimes \frac{\mathbf{M}\mathbf{H}^T}{(\mathbf{W}\mathbf{H})\mathbf{H}^T}$$



Towards Hybrid deep learning

... some prior works in Audio

- Coupling signal processing modules with deep learning for audio synthesis
- The example of DDSP (Engel & al.)



X. Wang & al. "Neural Source-Filter Waveform Models for Statistical Parametric Speech Synthesis," in IEEE/ACM Trans. on ASLP Proc., vol. 28, 2020.
J. Engel & al., "DDSP: Differentiable Digital Signal Processing," in Int. Conf. on Learning Representations (ICLR), 2020.

Towards Hybrid deep learning

... some prior works in Audio

- Phase retrieval from the magnitude spectrogram

$$\text{Find } \mathbf{X} \text{ s.t. } |X[\omega, \tau]| = A[\omega, \tau]$$

- The classic Griffin-Lim Algorithm (GLA)**

- Exploits spectrogram consistency
(\mathbf{X} should correspond to the complex spectrogram of a time domain signal x)

$$\text{Find } \mathbf{X} \text{ s.t. } \begin{cases} |X[\omega, \tau]| = A[\omega, \tau] \\ \mathbf{X} \in \text{Im}(\mathcal{G}) \end{cases}$$

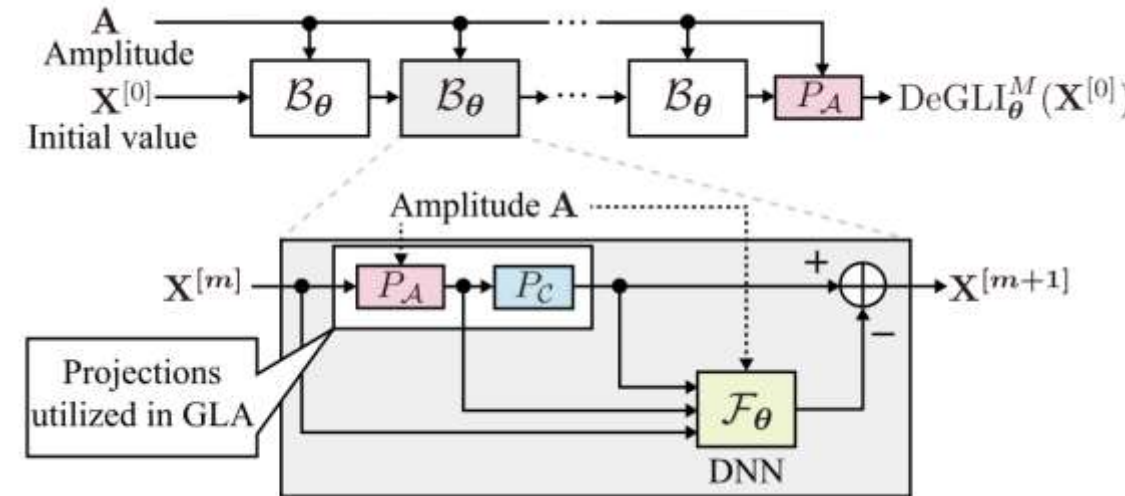
- Implemented as an iterative algorithm

$$\mathbf{X}^{[m+1]} = P_C(P_A(\mathbf{X}^{[m]}))$$

$$P_A(\mathbf{X})[\omega, \tau] = A[\omega, \tau] \frac{X[\omega, \tau]}{|X[\omega, \tau]|} \quad P_C(\mathbf{X}) = \mathcal{G}(\mathcal{G}^\dagger(\mathbf{X}))$$

- $\mathcal{G}, \mathcal{G}^\dagger$ are respectively the STFT and ISTFT operators

Deep griffin-Lim



Towards Hybrid deep learning

... some prior works in Audio

-And other very recent examples for virtual analog modelling combining Ordinary Differential Equations (ODEs) and neural network to learn the derivative function....
- ... at DAFx 2022 !!

Proceedings of the 25th International Conference on Digital Audio Effects (DAFx20in22), Vienna, Austria, September 6-10, 2022

VIRTUAL ANALOG MODELING OF DISTORTION CIRCUITS USING NEURAL ORDINARY DIFFERENTIAL EQUATIONS

*Jan Wilczek**

WolfSound
Katowice, Poland
jan.wilczek@thewolfsound.com

Alec Wright and Vesa Välimäki[†]

Acoustics Lab
Dept. Signal Processing and Acoustics
Aalto University, Espoo, Finland
alec.wright@aalto.fi

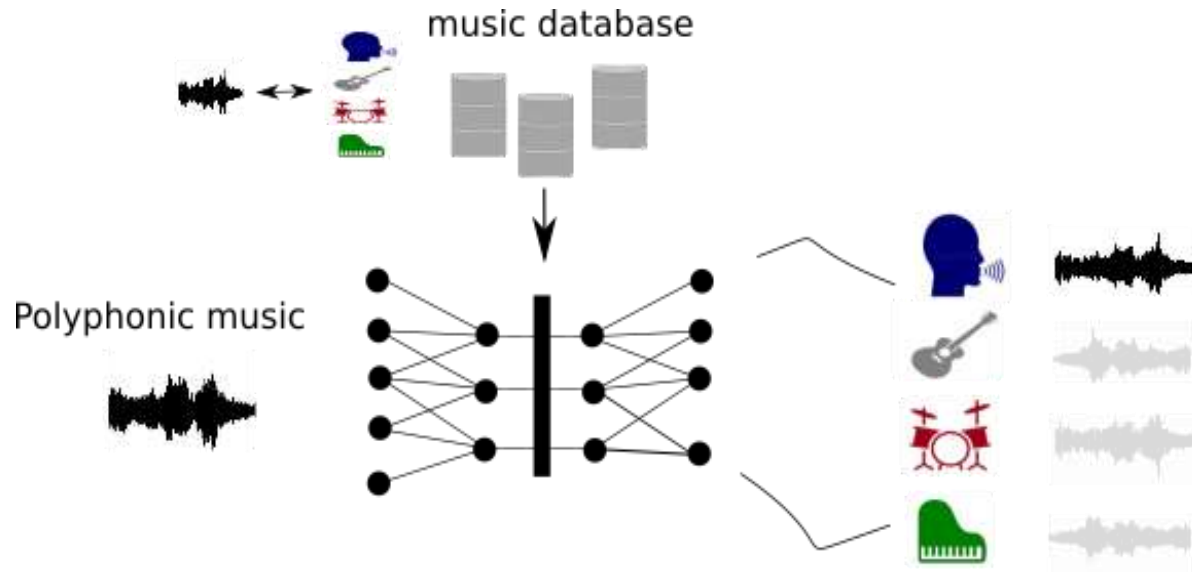
Emanuël A. P. Habets

International Audio Laboratories Erlangen[‡]
Erlangen, Germany
emanuel.habets
@audiolabs-erlangen.de

Towards Hybrid deep learning

... by **integrating our prior knowledge** about the nature of the processed data.

- For example in music source separation

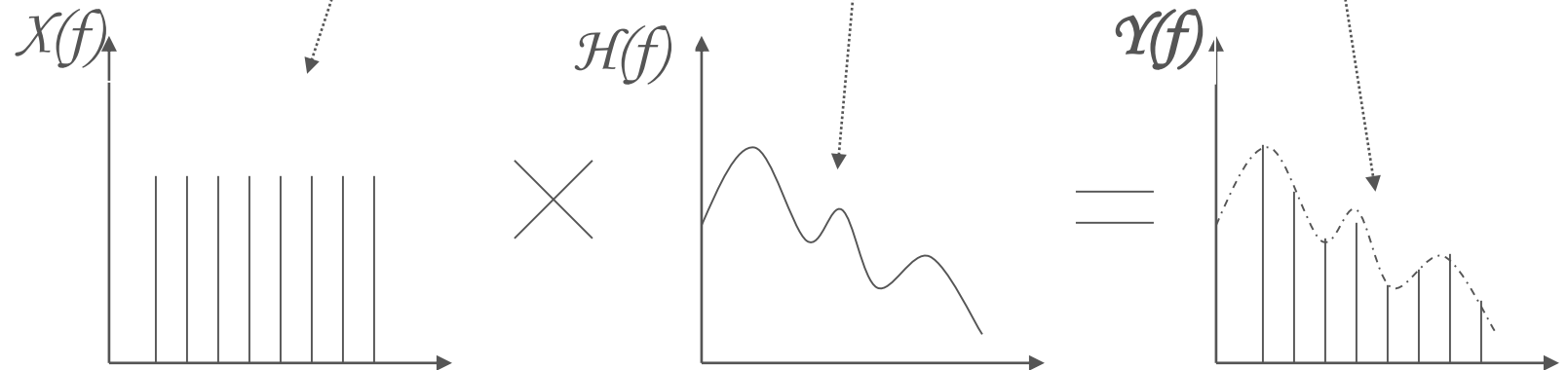


Main limitations:

- *Difficulty to obtain « aligned » data*
- *Knowledge learned (only) from data*
- *Complexity: overparametrized models*
- *Overconsumption regime*
- **Non-interpretable/non-controllable**

The source filter model

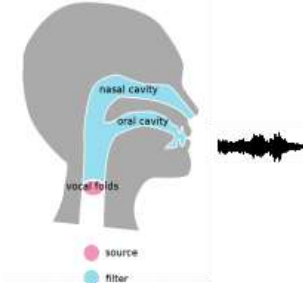
an efficient speech production model



Towards Hybrid deep learning

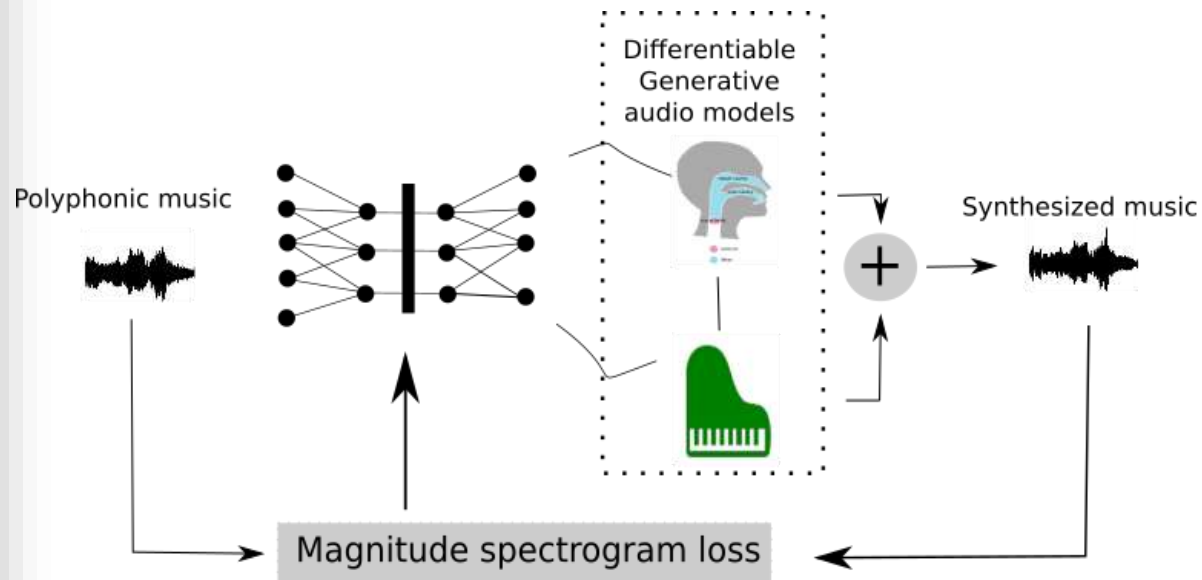
... by integrating our prior knowledge about the nature of the processed data.

Knowledge about « how the sound is produced » (e.g. sound production models)



Singing voice as a source / filter model :

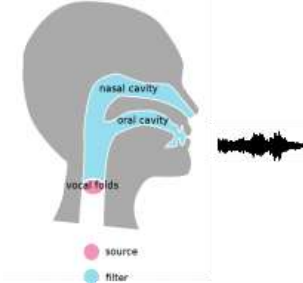
- source = vibration of vocal folds
- Filter = resonances of vocal/nasal cavities



Towards Hybrid deep learning

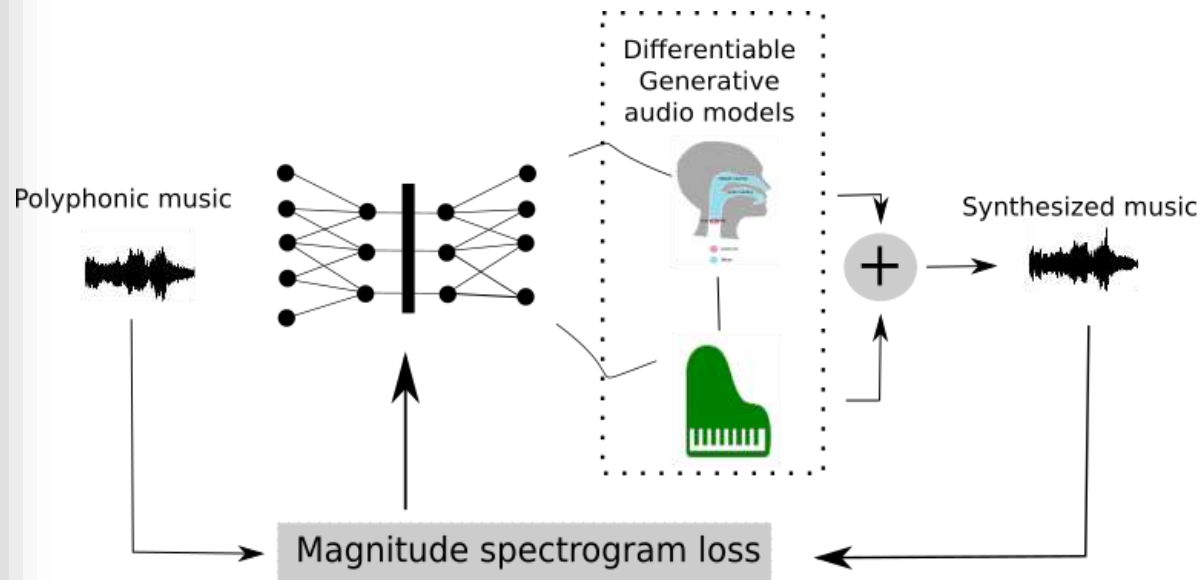
... by integrating our prior knowledge about the nature of the processed data.

Knowledge about « how the sound is produced » (e.g. sound production models)



Singing voice as a source / filter model :

- source = vibration of vocal folds
- Filter = resonances of vocal/nasal cavities



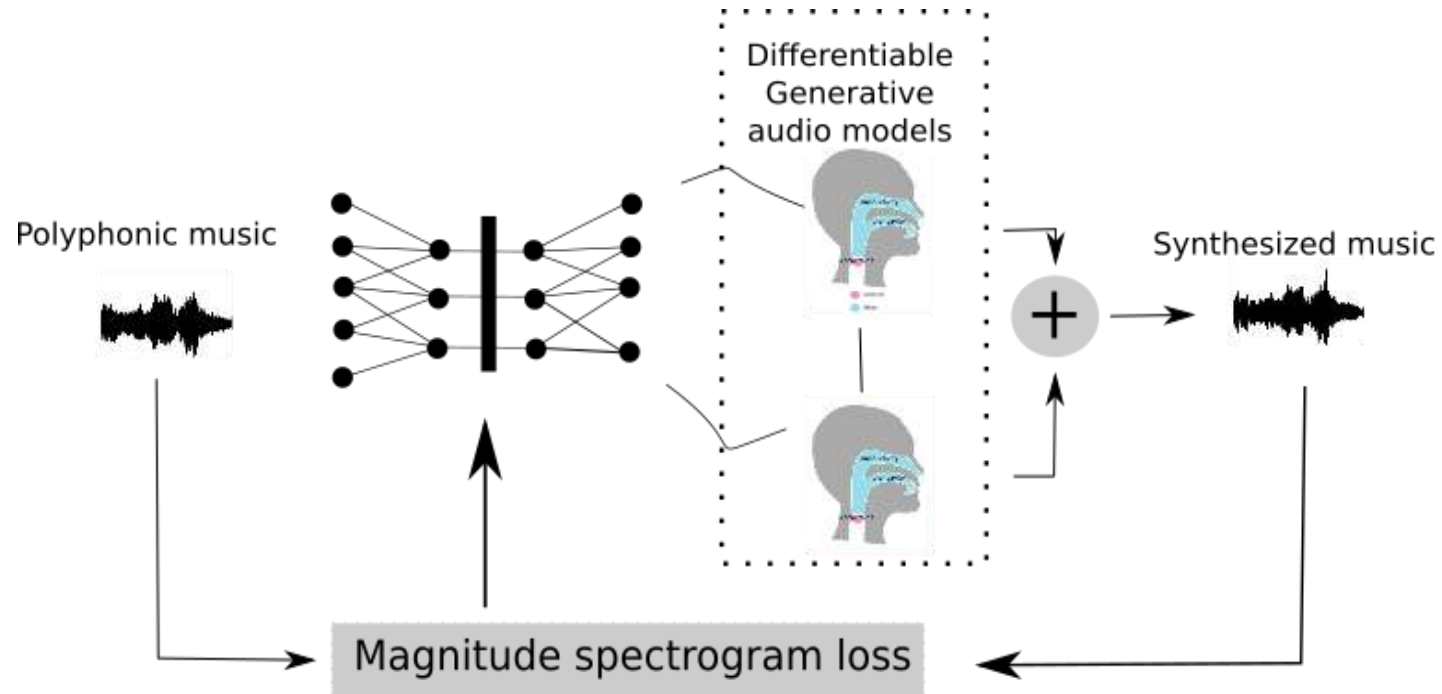
A new paradigm

- Model is at the « core » of neural architecture
- Source separation **by synthesis** (*no interference from other sources*)
- Learning only from the polyphonic recording (*no need of the true individual tracks*)

Towards Hybrid deep learning

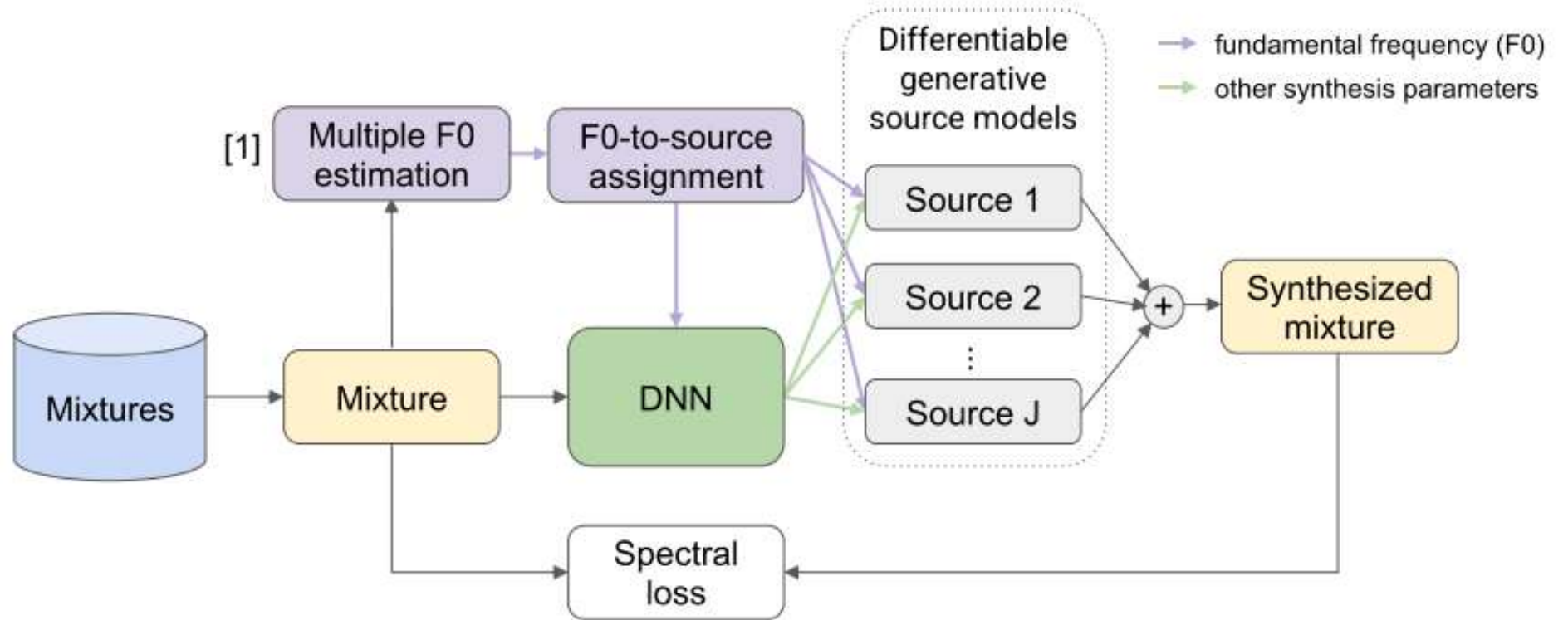
... by **integrating our prior knowledge** about the nature of the processed data.

- Preliminary work on source separation (choir singing)



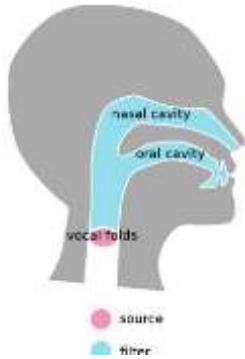
Unsupervised learning strategy

(e.g. no need of the individual source signals)

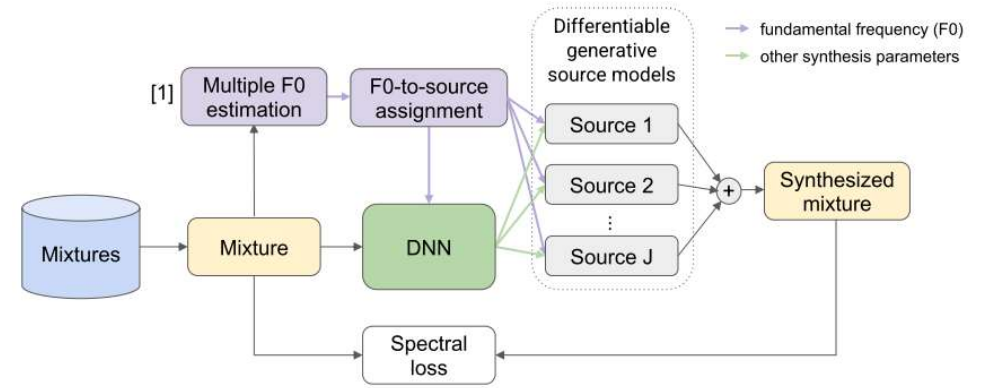
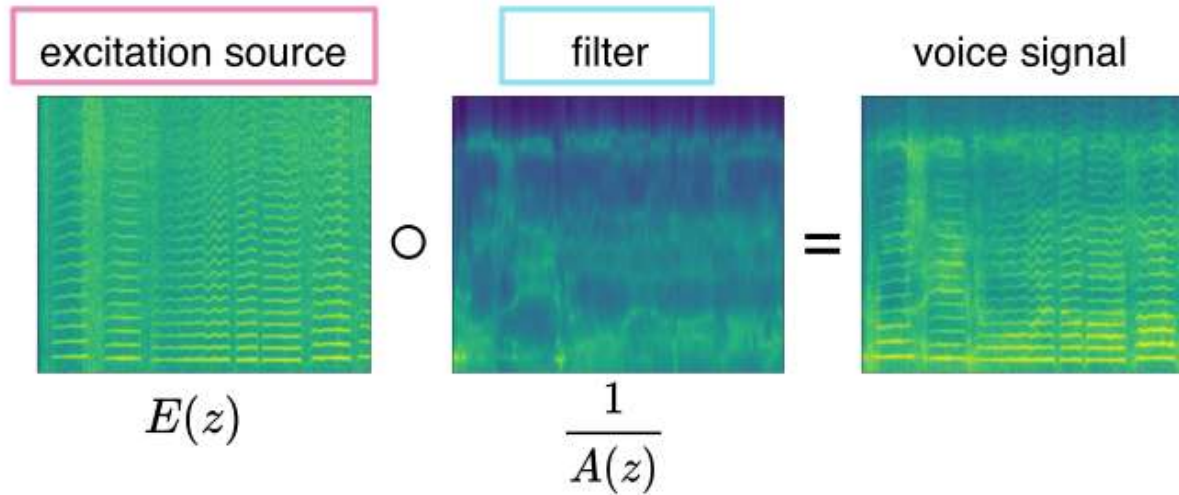


Parametric source models

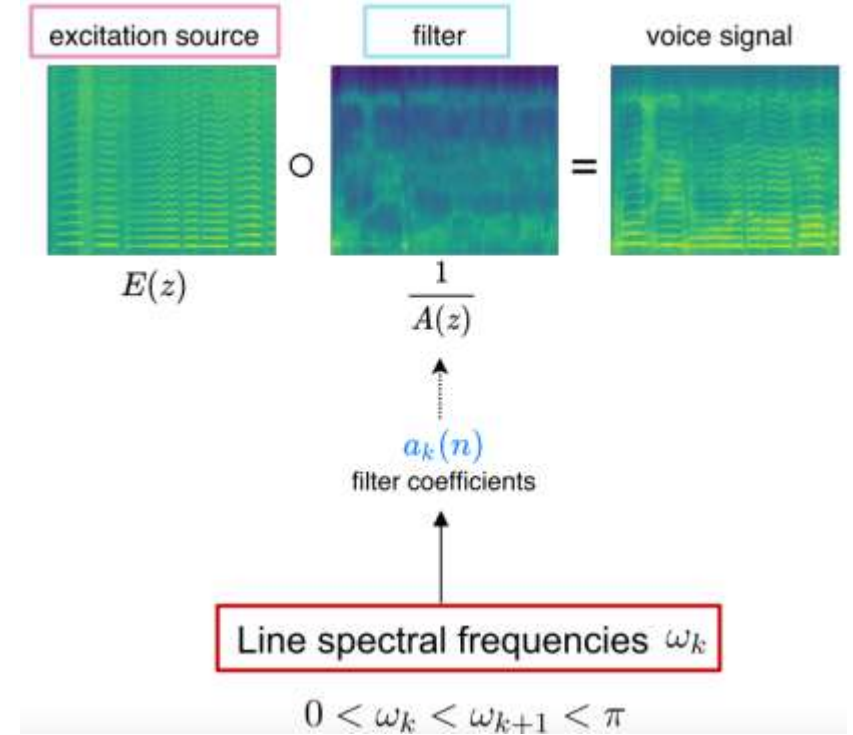
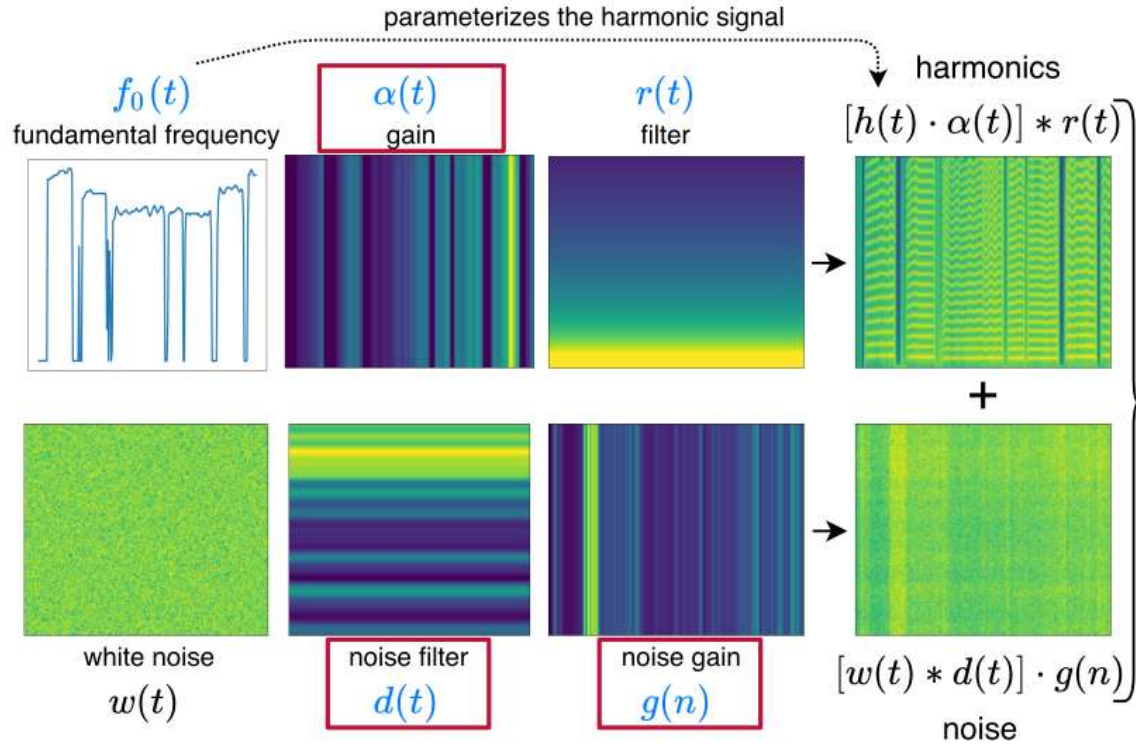
Singing voice as a source / filter model :



- source = vibration of vocal folds
- Filter = resonances of vocal/nasal cavities



Parametric source models



The Line Spectral frequencies (LSF)

- An alternative to the linear prediction coefficients

$$y(n) = x(n) + e(n) = x(n) - \sum_{i=1}^P a_i x(n-i)$$

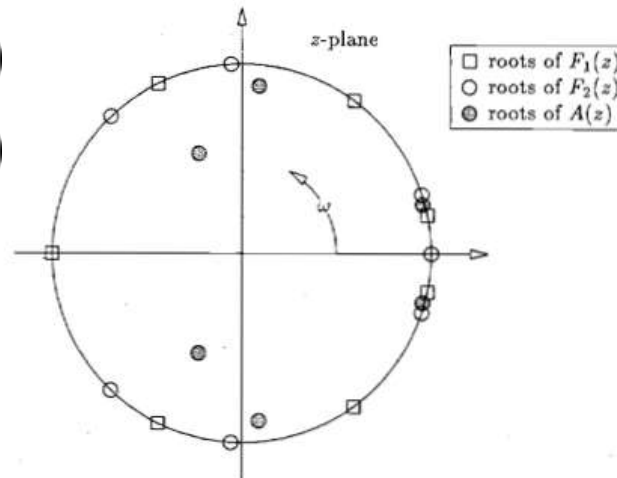
$$y(n) = x(n) \cdot (1 - \sum_{i=1}^P a_i z^{-i})$$

- Use of two auxilliary polynomials

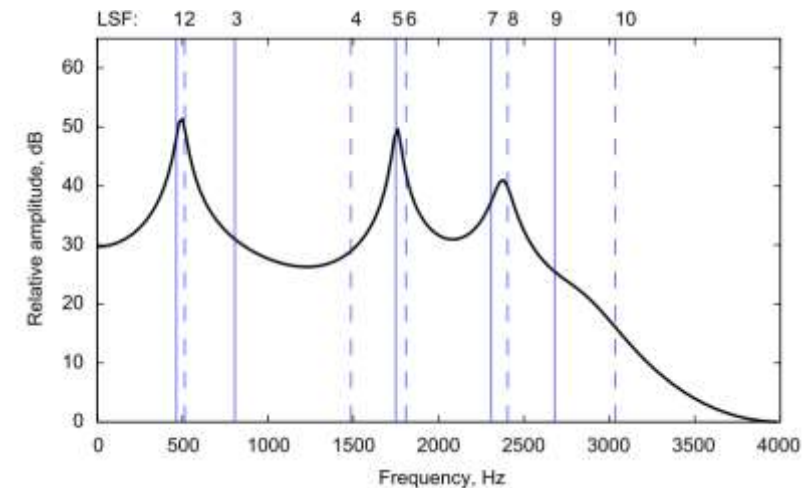
$$Y(z) = A(z)E(z) = (1 - a_1 z^{-1} - \dots - a_P z^{-P})E(z)$$

$$F_1(z) = A(z) + z^{-(P+1)} A(z^{-1})$$

$$F_2(z) = A(z) - z^{-(P+1)} A(z^{-1})$$



An example of LPC spectrum (from [2])



Unsupervised learning strategy

(e.g. no need of the individual source signals)

- A multi-scale spectral loss

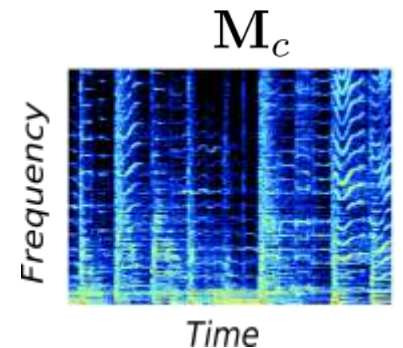
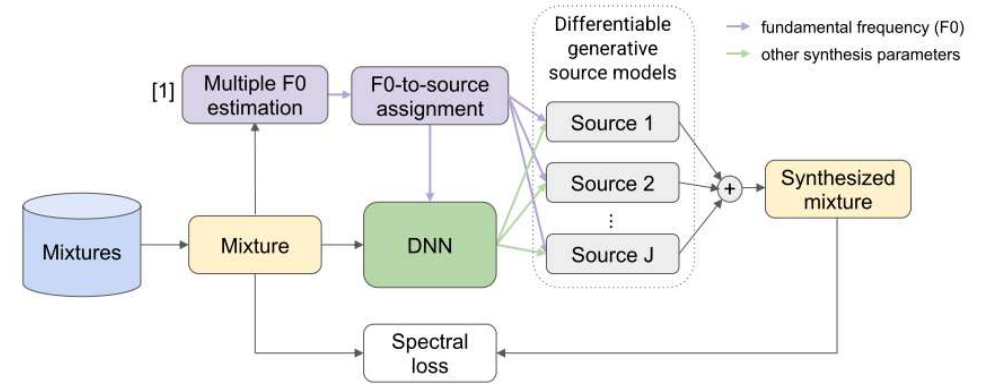
$$\mathcal{L}_{rec} = \sum_c \mathcal{L}_c$$

With

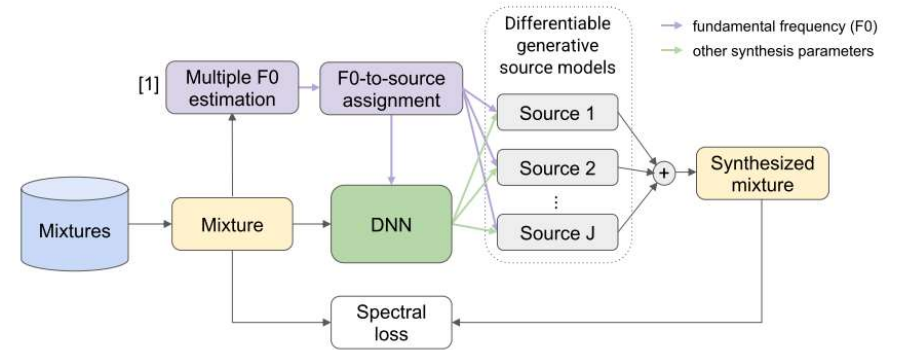
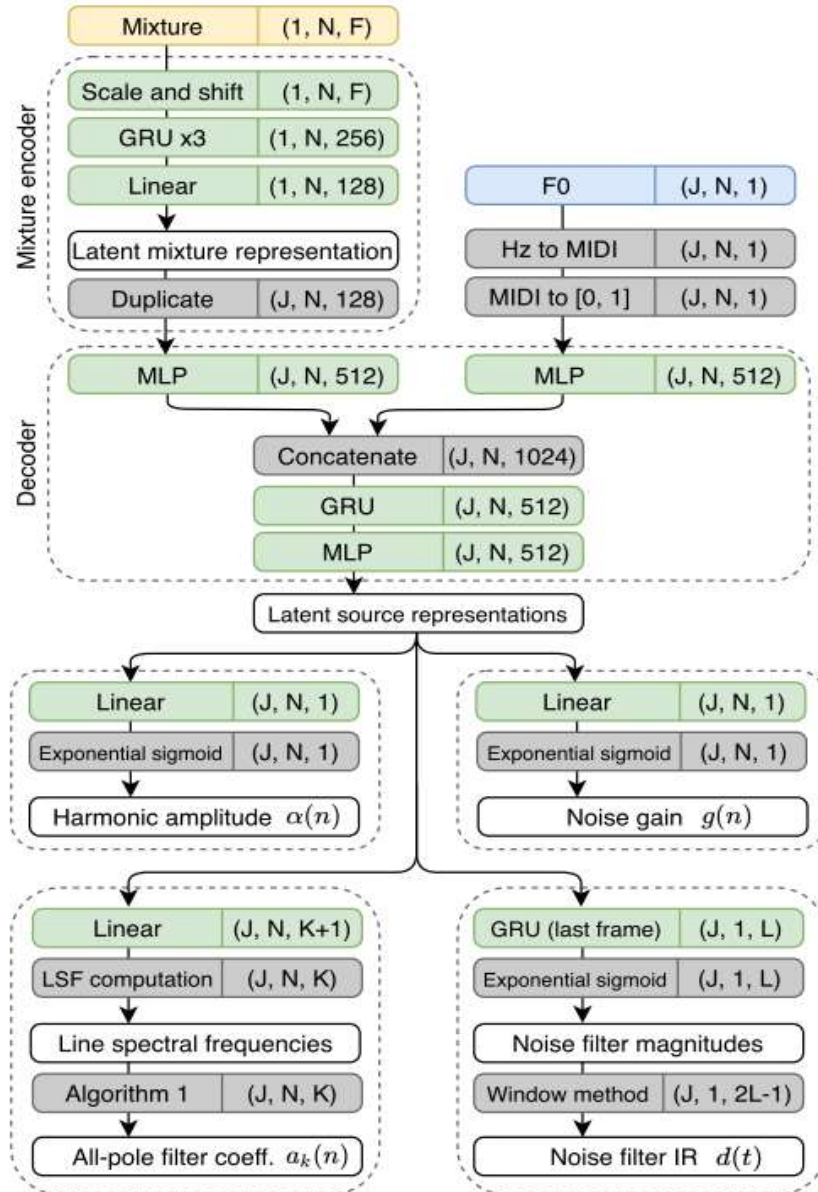
$$\mathcal{L}_c = \|\mathbf{M}_c - \tilde{\mathbf{M}}_c\|_1 + \|\log(\mathbf{M}_c) - \log(\tilde{\mathbf{M}}_c)\|_1$$

Where \mathbf{M}_c and $\tilde{\mathbf{M}}_c$ denote the magnitude spectrograms of the input mixture and its estimate, respectively,

and $c = [2048, 1024, 512, 256, 128, 64]$ indicates the FFT size used to compute the STFT. The frames overlap by 75%.

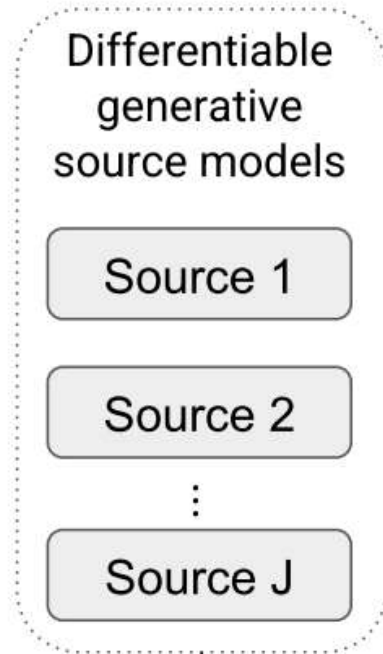


Global architecture overview

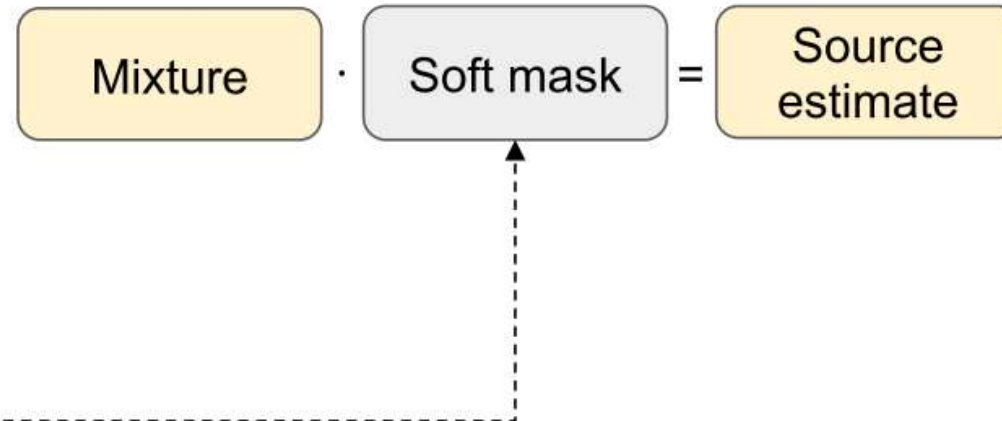


Synthesis or filtering

Synthesis

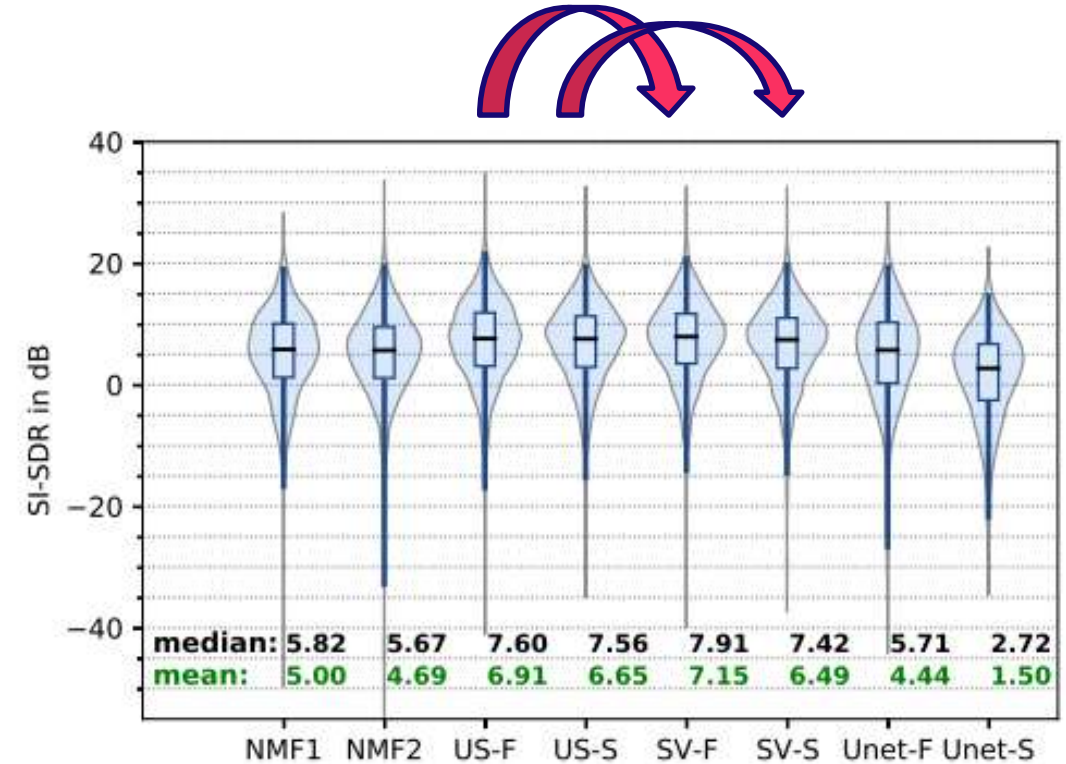


Filtering



Some results

- Unsupervised (US) \approx supervised (SU)



(b) $J = 4$ sources



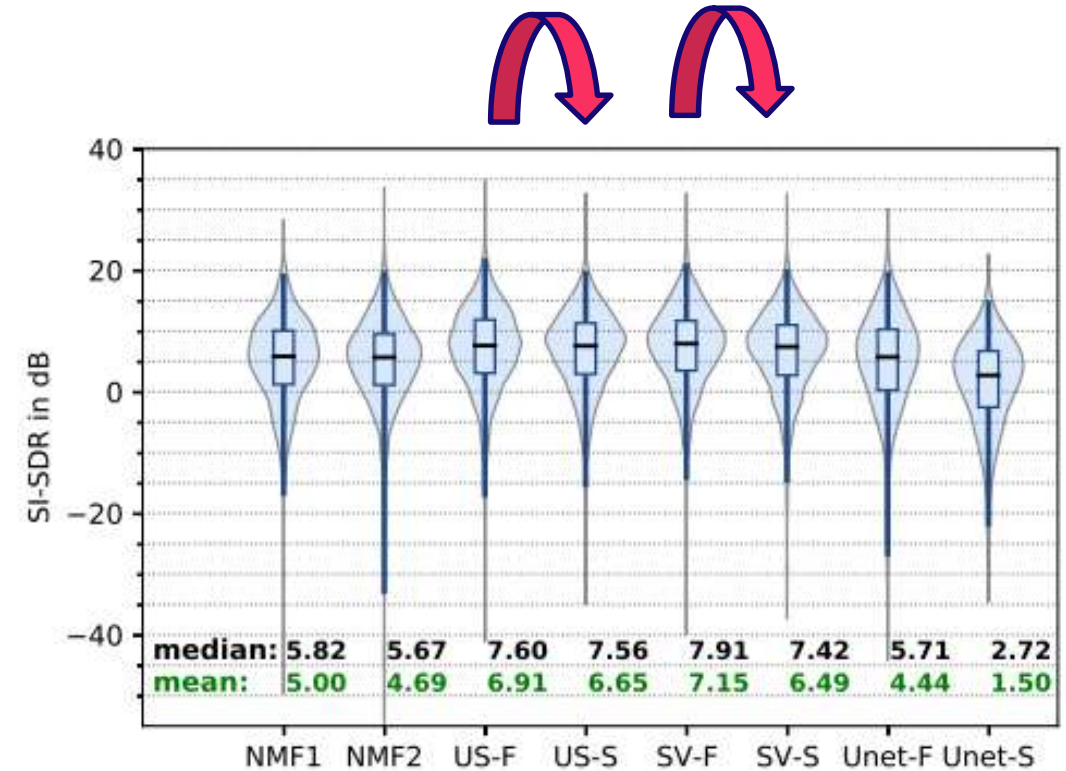
NMF1: S. Ewert and M. Müller, "Using score-informed constraints for NMF-based source separation," in Proc. IEEE Int. Conf. on Acoustics, Speech and Signal Processing. IEEE, 2012, pp. 129–132.

NMF2: J.-L. Durrieu, B. David, and G. Richard, "A musically motivated mid-level representation for pitch estimation and musical audio source separation," IEEE J. Selected Topics in Signal Processing, vol. 5, no. 6, pp. 1180–1191, 2011.

UNET: D. Petermann, P. Chandna, H. Cuesta, J. Bonada, and E. Gomez, "Deep learning based source separation applied to choir ensembles," in Proc. Int. Soc. Music Inf. Retrieval Conf., 2020, pp. 733–739.

Some results

- Unsupervised (US) \approx supervised (SU)
- Almost no drop of performances when using only 3% of the training data (US-F vs. US-S and SV-F vs. SV-S)



(b) $J = 4$ sources



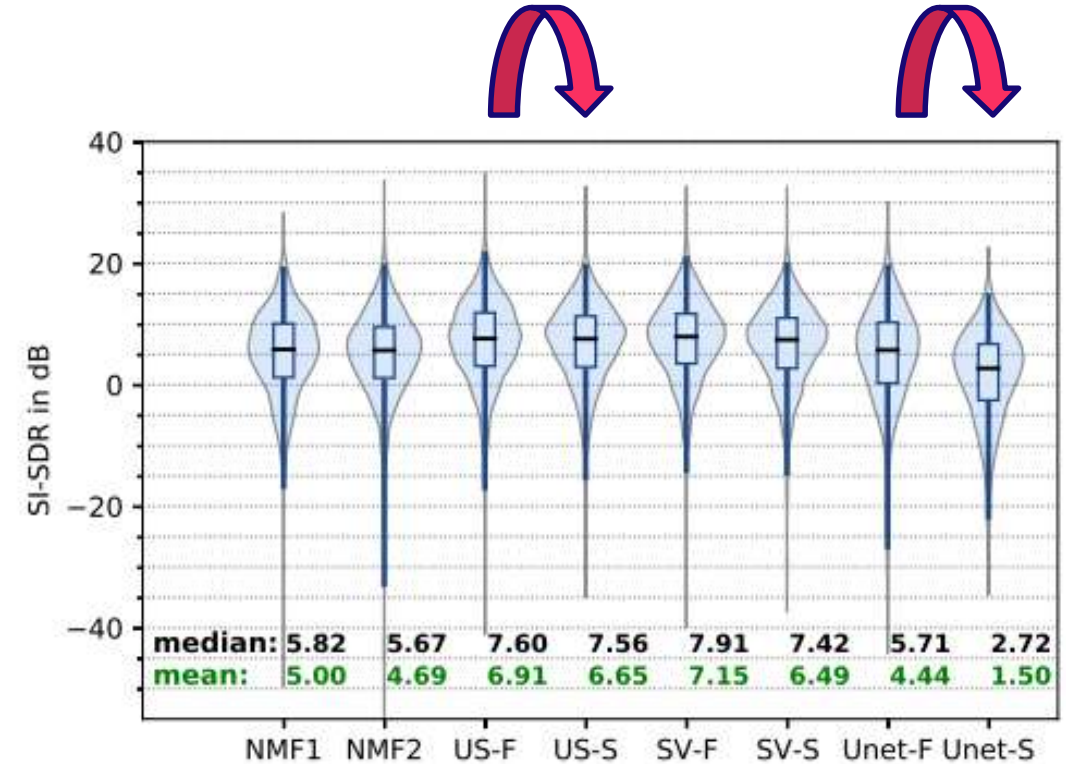
NMF1: S. Ewert and M. Müller, "Using score-informed constraints for NMF-based source separation," in Proc. IEEE Int. Conf. on Acoustics, Speech and Signal Processing. IEEE, 2012, pp. 129–132.

NMF2: J.-L. Durrieu, B. David, and G. Richard, "A musically motivated mid-level representation for pitch estimation and musical audio source separation," IEEE J. Selected Topics in Signal Processing, vol. 5, no. 6, pp. 1180–1191, 2011.

UNET: D. Petermann, P. Chandna, H. Cuesta, J. Bonada, and E. Gomez, "Deep learning based source separation applied to choir ensembles," in Proc. Int. Soc. Music Inf. Retrieval Conf., 2020, pp. 733–739.

Some results

- Unsupervised (US) \approx supervised (SU)
- Almost no drop of performances when using only 3% of the training data (US-F vs. US-S and SV-F vs. SV-S)
- ..much larger drop of performances of the supervised baseline model (Unet)



(b) $J = 4$ sources



NMF1: S. Ewert and M. Mueller, "Using score-informed constraints for NMF-based source separation," in Proc. IEEE Int. Conf. on Acoustics, Speech and Signal Processing. IEEE, 2012, pp. 129–132.

NMF2: J.-L. Durrieu, B. David, and G. Richard, "A musically motivated mid-level representation for pitch estimation and musical audio source separation," IEEE J. Selected Topics in Signal Processing, vol. 5, no. 6, pp. 1180–1191, 2011.

UNET: D. Petermann, P. Chandna, H. Cuesta, J. Bonada, and E. Gomez, "Deep learning based source separation applied to choir ensembles," in Proc. Int. Soc. Music Inf. Retrieval Conf., 2020, pp. 733–739.

A short audio demo and some take aways

- **A short demo at**
 - <https://schufo.github.io/umss/>
 - Ou [lien local](#)
- **Some take aways**
 - Only a small amount of data needed
 - Filtering the mixture better than synthesis
 - Differentiable stable all-pole filter
 - Parameterization of the mixture is provided



To conclude

- The potential for hybrid deep learning ...
 - **Interpretability, Controllability, Explainability**
 - Hybrid model becomes controllable by human-understandable parameters
 - New audio capabilities: perceptually meaningful sound transformation
 - **Frugality: gain of several orders of magnitude** in the need of data and model complexity
 - **Towards a more resource efficient and sustainable AI**