

MOTIVATION

- Nonnegative matrix factorization (NMF) can be used to decompose a spectrogram $\mathbf{V} \in \mathbb{R}^{M \times N}$ into two nonnegative latent factors $\mathbf{W} \in \mathbb{R}^{M \times K}$ and $\mathbf{H} \in \mathbb{R}^{K \times N}$ which respectively encode spectral patterns (dictionary) and how these are mixed (activation).
- Results depend heavily on the time-frequency transform used for computing \mathbf{V} .
- Can we learn a transform Φ so that \mathbf{V} can be well approximated using NMF?

BASELINE : IS-NMF

Audio data

$\mathbf{Y} \in \mathbb{R}^{M \times N}$: matrix that contains N adjacent and overlapping short-time M -wide frames of the sound sample y

IS-NMF with sparsity

Minimize

$$D(|\Phi_{\text{DCT}} \mathbf{Y}|^{\circ 2} | \mathbf{W} \mathbf{H}) + \lambda \frac{M}{K} \|\mathbf{H}\|_1$$

s.t. $\mathbf{W} \geq 0, \mathbf{H} \geq 0, \forall k, \|\mathbf{w}_k\|_1 = 1$ (1)

with $D(\mathbf{A}|\mathbf{B}) = \sum_{ij} (a_{ij}/b_{ij} - \log(a_{ij}/b_{ij}) - 1)$ (Itakura-Saito divergence), factorization rank K

TRANSFORM LEARNING

Proposed TL-NMF (inspired from [1])

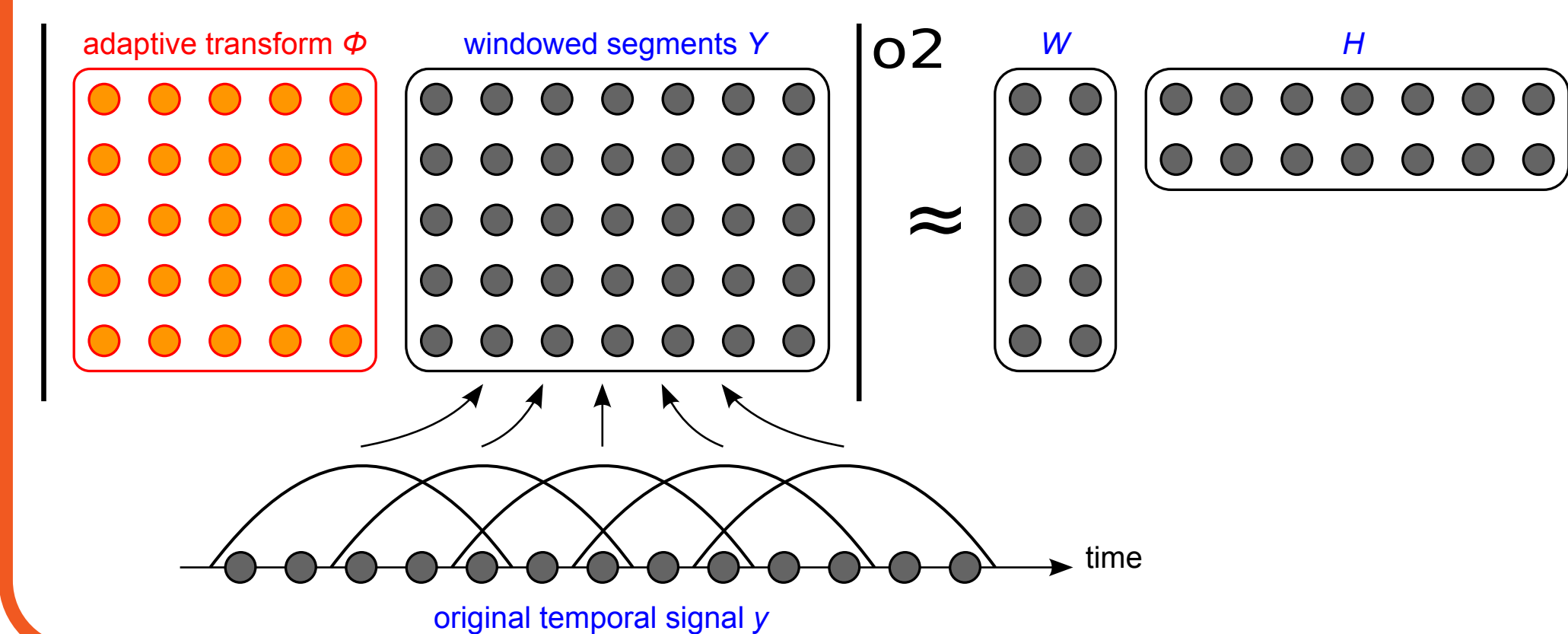
Minimize

$$C_\lambda(\Phi, \mathbf{W}, \mathbf{H}) \stackrel{\text{def}}{=} D(|\Phi \mathbf{Y}|^{\circ 2} | \mathbf{W} \mathbf{H}) + \lambda \frac{M}{K} \|\mathbf{H}\|_1$$

s.t. $\mathbf{W} \geq 0, \mathbf{H} \geq 0, \forall k, \|\mathbf{w}_k\|_1 = 1, \Phi^T \Phi = \mathbf{I}_M$ (2)

Orthogonal constraint on Φ

- Gently departs from Φ_{DCT}
- Avoids singularity along with trivial solutions such as $(\Phi, \mathbf{W}, \mathbf{H}) = (\mathbf{0}, \mathbf{0}, \mathbf{0})$
- Easy inversion for synthesis



PROPOSED ALGORITHM

Algorithm 1: TL-NMF

Input : $\mathbf{Y}, \tau, K, \lambda$

Output: $\Phi, \mathbf{W}, \mathbf{H}$

Initialize Φ, \mathbf{W} and \mathbf{H} at random

while $\epsilon > \tau$ **do**

 Update \mathbf{H}

 Update \mathbf{W}

 Update Φ (new)

 Compute stopping criterion ϵ

end

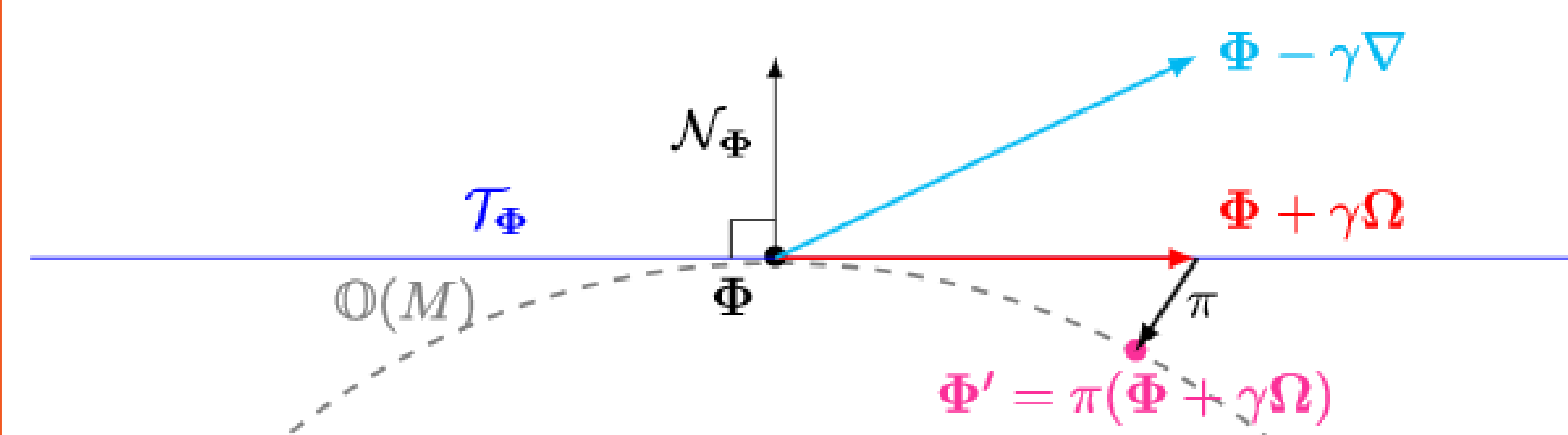
Update of \mathbf{H} and \mathbf{W}

Majoration-minimization leading to standard multiplicative updates [2]

Update of Φ

Projected gradient descent onto the orthogonal matrices manifold following [3]

- 1) Compute gradient ∇ of the objective function
- 2) Compute natural gradient $\Omega = \Phi \nabla^T \Phi - \nabla$
- 3) Find a suitable stepsize γ satisfying Armijo rule on the manifold
- 4) Update the transform via a projection onto the manifold as $\Phi \leftarrow \pi(\Phi + \gamma \Omega)$
- 5) Resolve sign ambiguity on Φ by imposing its first column entries to be positive



REFERENCES

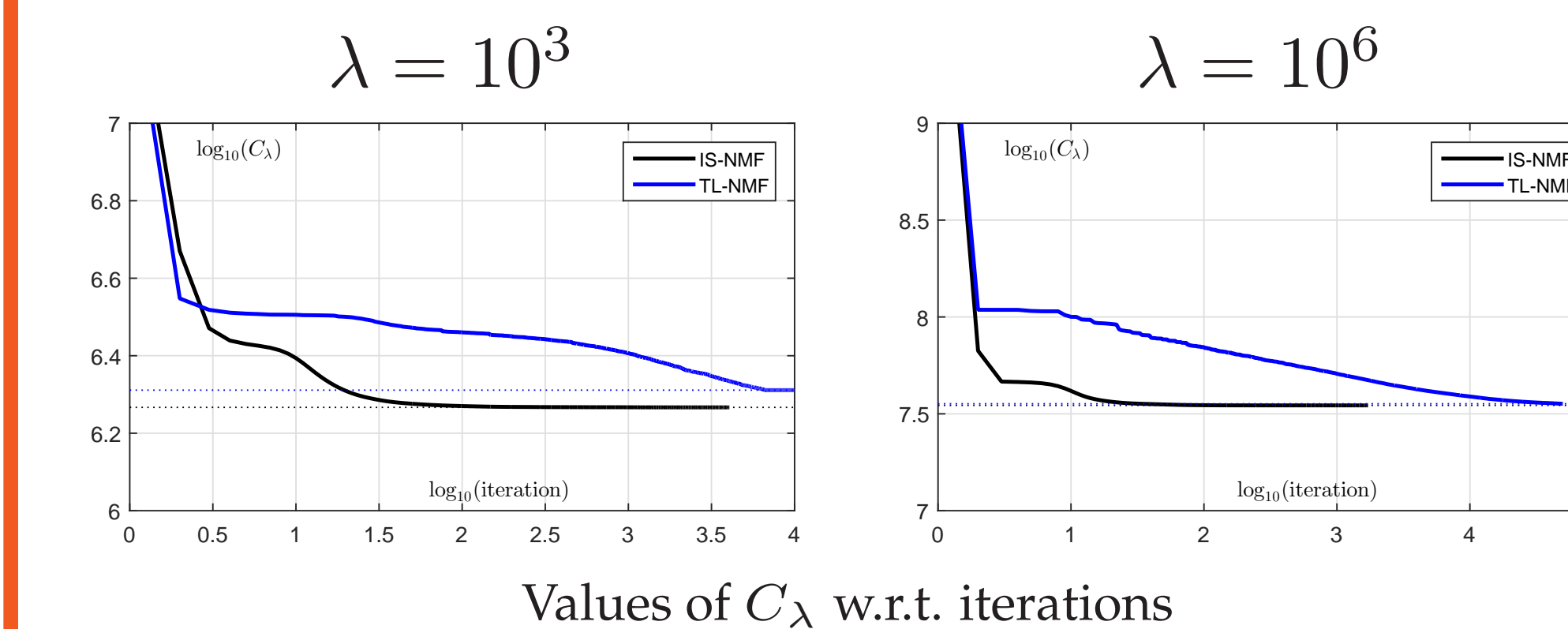
- [1] S. Ravishanker and Y. Bresler, "Learning sparsifying transforms," *IEEE T. Signal Process.*, 2013.
- [2] C. Févotte and J. Idier, "Algorithms for nonnegative matrix factorization with the β -divergence," *Neural Comput.*, 2011.
- [3] J. H. Manton, "Optimization algorithms exploiting unitary constraints," *IEEE T. Signal Process.*, 2002.

MUSIC DECOMPOSITION

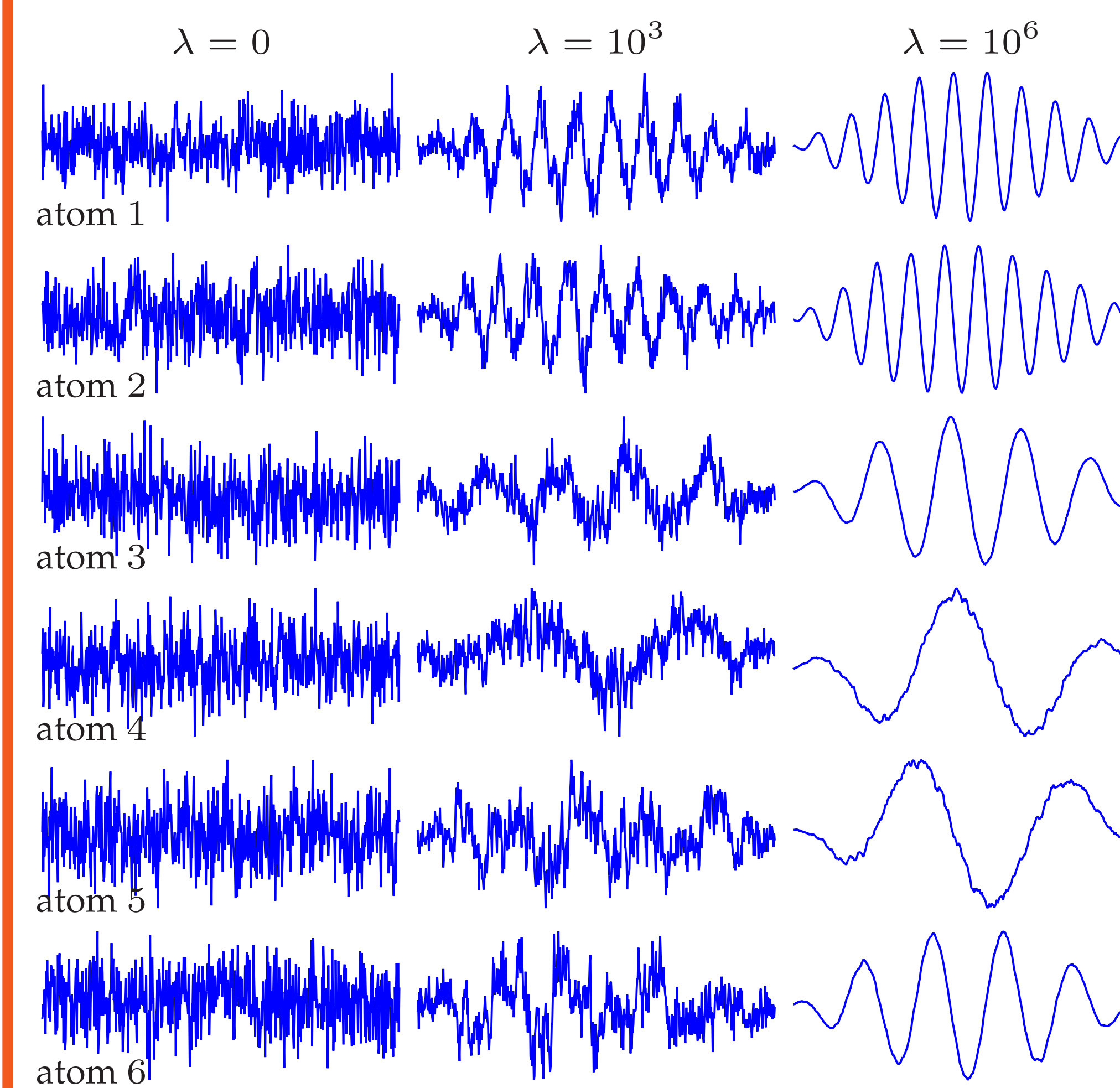
Setup

23 s long excerpt from *Mamavatu* by Susheela Raman using 50% overlapping 40 ms-long sine bell windows with factorization rank $K = 10$

Results

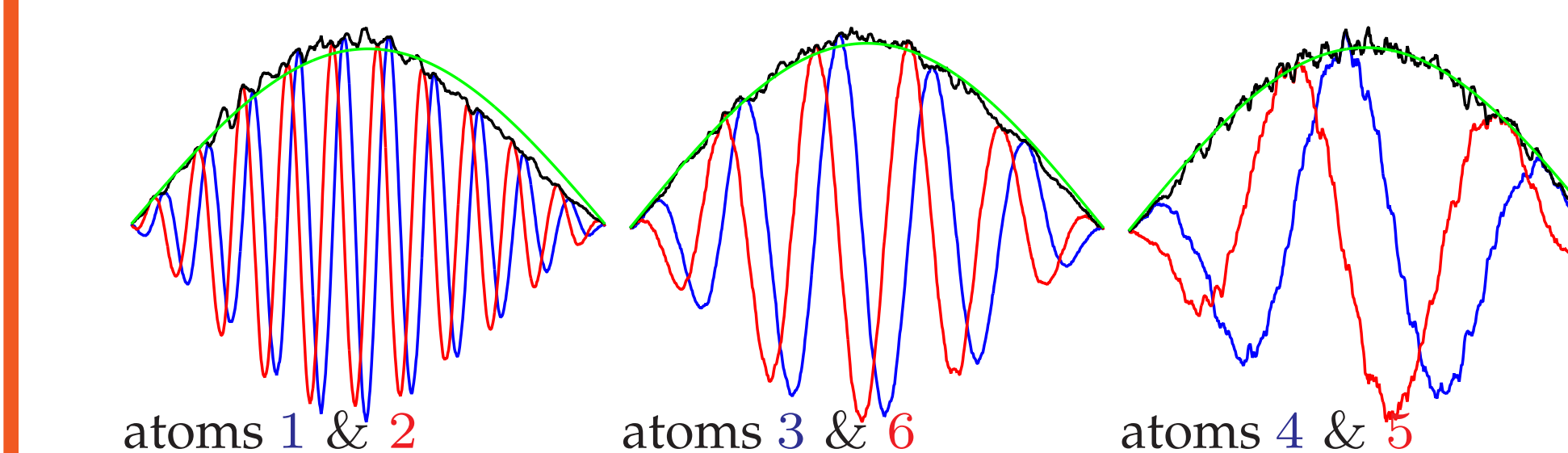


TL-NMF reaches similar objective function values despite random initialization



Six most significant atoms learnt by TL-NMF from random initializations

Rows of Φ become oscillatory and smoother as λ increases



Atoms form pairs in phase quadrature

SUPERVISED SEPARATION

Setup

Separate a sound sample as $y = \hat{y}_{\text{sp}} + \hat{y}_{\text{no}}$ based on the reference data

- $\mathbf{Y}_{\text{sp}} \in \mathbb{R}^{M \times N_{\text{sp}}}$, speech female speaker (21 s)
- $\mathbf{Y}_{\text{no}} \in \mathbb{R}^{M \times N_{\text{no}}}$, bus noise (30 s)

as

$$\mathbf{V} \approx \mathbf{W}_{\text{sp}} \mathbf{H}_{\text{sp}} + \mathbf{W}_{\text{no}} \mathbf{H}_{\text{no}} \quad (3)$$

where $\mathbf{W}_{\text{sp}} = |\Phi_{\text{DCT}} \mathbf{Y}_{\text{sp}}|^{\circ 2}$, $\mathbf{W}_{\text{no}} = |\Phi_{\text{DCT}} \mathbf{Y}_{\text{no}}|^{\circ 2}$ and $\mathbf{H}_{\text{sp}}, \mathbf{H}_{\text{no}}$ are subject to a sparsity constraint.

Separation with TL-NMF

Minimize

$$C_\lambda(\Phi, \mathbf{H}_{\text{sp}}, \mathbf{H}_{\text{no}}) \stackrel{\text{def}}{=} \quad (4)$$

$$D(|\Phi \mathbf{Y}|^{\circ 2} | |\Phi \mathbf{Y}_{\text{sp}}|^{\circ 2} \mathbf{H}_{\text{sp}} + |\Phi \mathbf{Y}_{\text{no}}|^{\circ 2} \mathbf{H}_{\text{no}})$$

$$+ \lambda_{\text{sp}} \frac{M}{N_{\text{sp}}} \|\mathbf{H}_{\text{sp}}\|_1 + \lambda_{\text{no}} \frac{M}{N_{\text{no}}} \|\mathbf{H}_{\text{no}}\|_1$$

$$\text{s.t. } \Phi^T \Phi = \mathbf{I}, \mathbf{H}_{\text{sp}} \geq 0, \mathbf{H}_{\text{no}} \geq 0,$$

Φ now appears in both sides of the divergence

Results

Sound sample generated by mixing a speech utterance with a bus noise at two different SNR

Comparison using BSS_eval metrics with baseline ($\hat{y}_{\text{sp}} = \hat{y}_{\text{no}} = y/2$) and IS-NMF with sparsity

$\lambda_{\text{sp}} = \lambda_{\text{no}} = \lambda$ was fixed manually

Method	SDR (dB)		SIR (dB)		SAR (dB)	
SNR = -10 dB	\hat{y}_{sp}	\hat{y}_{no}	\hat{y}_{sp}	\hat{y}_{no}	\hat{y}_{sp}	\hat{y}_{no}
	Baseline	-9.50 10.00	-9.50 10.00	∞ ∞		
	IS-NMF	-6.75 6.82	-5.00 13.95	4.12 7.93		
TL-NMF	1.73 12.29	13.44 13.33	2.22 19.20			
SNR = 0 dB	\hat{y}_{sp}	\hat{y}_{no}	\hat{y}_{sp}	\hat{y}_{no}	\hat{y}_{sp}	\hat{y}_{no}
	Baseline	0.10 0.08	0.10 0.08	∞ ∞		
	IS-NMF	1.73 0.69	3.06 5.32	9.30 3.65		
TL-NMF	6.50 5.81	12.11 9.16	8.16 9.00			

CONCLUSION

- Introduction of transform learning for NMF
- Proposal of a new block-coordinate descent algorithm
- TL-NMF automatically uncovers data-driven oscillatory atoms