

Nonnegative Tucker Decomposition with Beta-divergence for Music Structure Analysis of Audio Signals

A. Marmoret¹, F. Voorwinden¹, V. Leplat², J.E. Cohen³, F. Bimbot¹

¹Univ Rennes, Inria, CNRS, IRISA, France, ²Center for Artificial Intelligence Technology (CAIT), Skoltech, Moscow, Russia, ³CNRS, CREATIS, Villeurbanne France - axe1.marmoret@irisa.fr

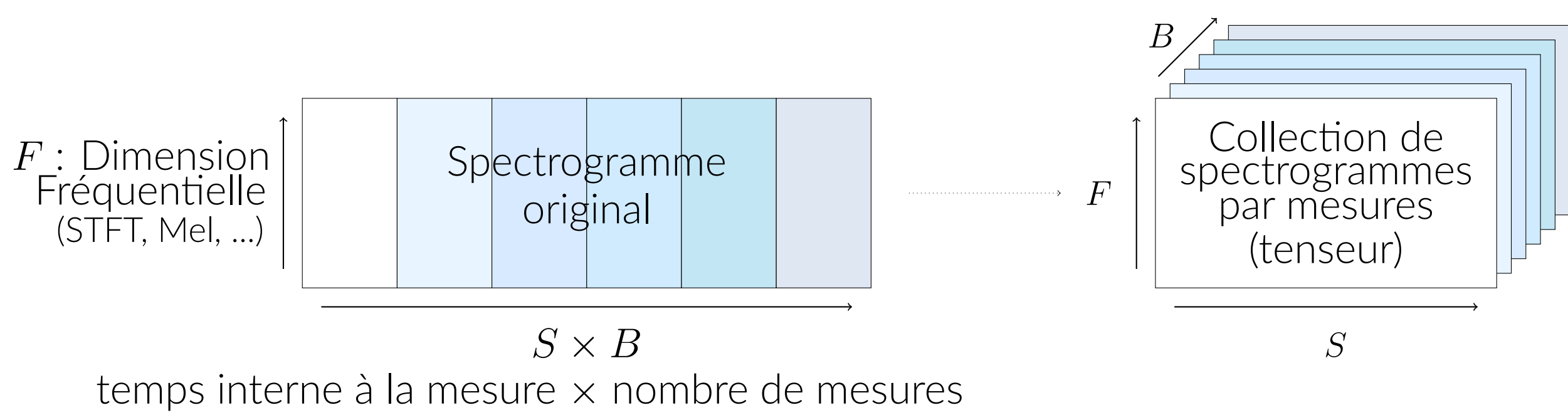
Résumé du poster

Ce poster présente une technique de factorisation appelée décomposition en Tucker nonnégatif (NTD), et un algorithme pour la calculer, optimisée selon la β -divergence [1].

La NTD permet d'extraire des patterns audios dans une musique, à l'échelle de la mesure [2]. Ces patterns peuvent être utilisés pour l'analyse musicale et la recomposition, ainsi que comme une représentation de mi-niveau, permettant notamment d'estimer la segmentation structurale.

Dans le cadre de l'analyse de la musique, la β -divergence est montrée plus pertinente que la norme euclidienne utilisée dans un travail antérieur [1].

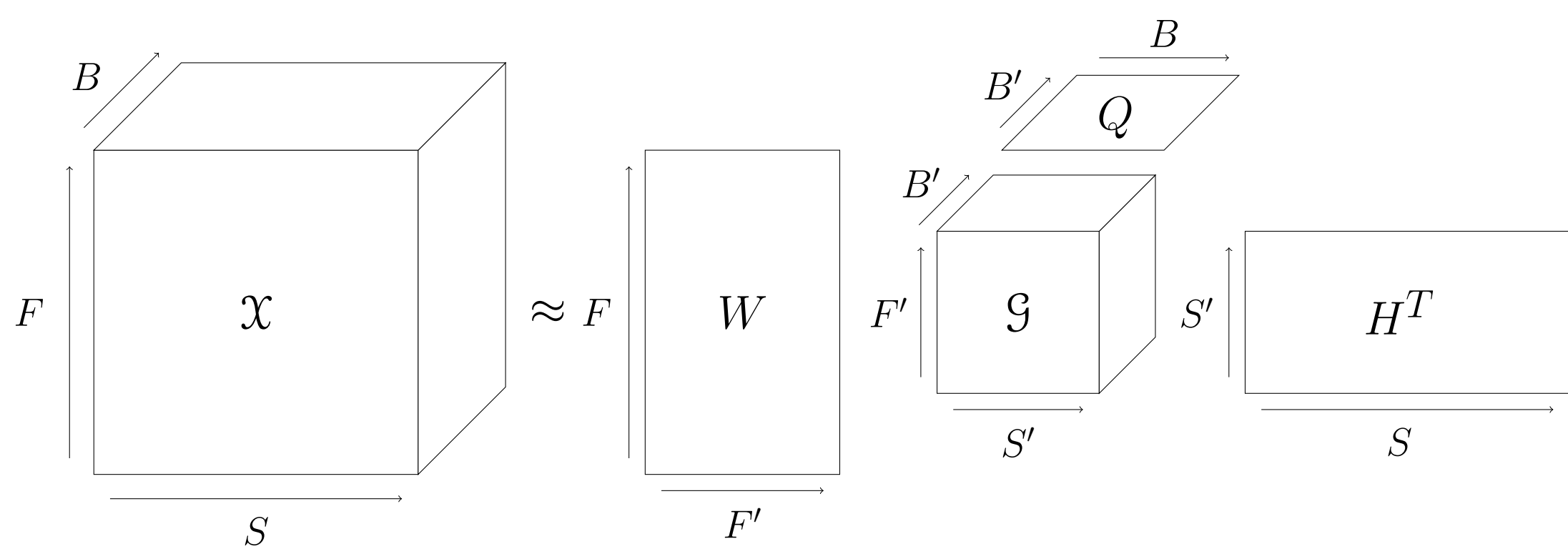
Etude à l'échelle de la mesure



Décomposition en Tucker nonnégatif (NTD)

Modèle NTD théorique :

$$\text{NTD} : \mathcal{X} \approx \mathcal{G} \times_1 W \times_2 H \times_3 Q$$



Pour chaque élément: $\mathcal{X}(f, s, b) \approx \sum_{f', s', b'=1}^{F', S', B'} \mathcal{G}(f', s', b') W(f, f') H(s, s') Q(b, b')$

NTD en pratique :

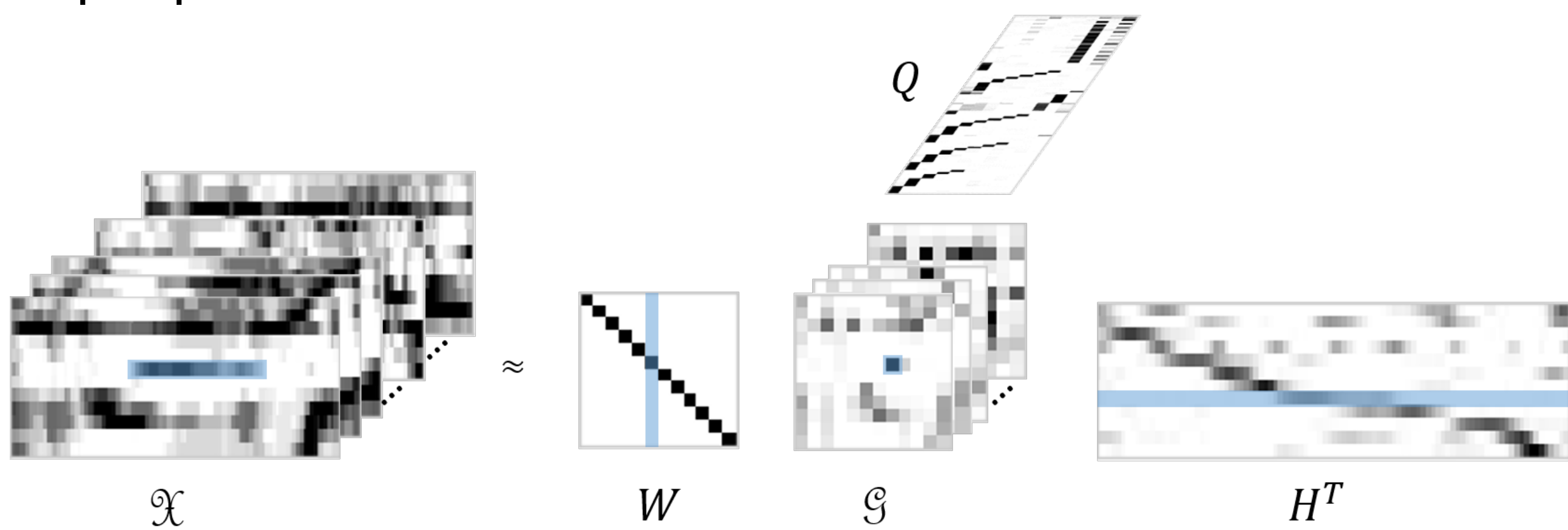


Figure. NTD sur le chromagramme de "Come Together" ($F' = 12, S' = 12$ et $B' = 10$).

β -divergence

$$d_{\beta}(x|y) = \begin{cases} \frac{x}{y} - \log\left(\frac{x}{y}\right) - 1 & \beta = 0 \text{ Itakura-Saito- (IS-) divergence} \\ x \log\left(\frac{x}{y}\right) + (y - x) & \beta = 1 \text{ Kullback-Leibler- (KL-) divergence} \\ \frac{x^{\beta} + (\beta - 1)y^{\beta} - \beta xy^{\beta - 1}}{\beta(\beta - 1)} & \beta \in \mathbb{R} \setminus \{0, 1\} \end{cases} \quad (1)$$

NTD : problème d'optimisation

• Problème d'optimisation β -NTD (d_{β} est la β -divergence relative à chaque élément) :

$$\arg \min_{W \geq 0, H \geq 0, Q \geq 0, \mathcal{G} \geq 0} d_{\beta}(\mathcal{X} | \mathcal{G} \times_1 W \times_2 H \times_3 Q)$$

L'algorithme [1] utilise les règles MU [3], revisités et optimisés pour l'algèbre tensoriel.

Algorithm 1: Une boucle de β -NTD

Input: $\mathcal{X}, \mathcal{G}, W, H, Q, \epsilon, \beta$

Output: \mathcal{G}, W, H, Q

$$V = (\mathcal{G} \times_2 H \times_3 Q)_{(1)}$$

$$W \leftarrow \max \left(W \cdot \left(\frac{[(WV)^{-(\beta-2)} \mathcal{X}_{(1)}] V^T}{(WV)^{-(\beta-1)} V^T} \right)^{\gamma(\beta)}, \epsilon \right) \quad (\text{Analogie pour } H \text{ et } Q)$$

/* Cette ligne utilise la propriété $\mathcal{G}_{(1)}(H \otimes Q)^T = (\mathcal{G} \times_2 H \times_3 Q)_{(1)}$, qui permet de réduire fortement la complexité de l'algorithme */

$$\mathcal{N} = (\mathcal{G} \times_1 W \times_2 H \times_3 Q)^{(\beta-2)} \cdot \mathcal{X}$$

$$\mathcal{D} = (\mathcal{G} \times_1 W \times_2 H \times_3 Q)^{(\beta-1)}$$

$$\mathcal{G} \leftarrow \max \left(\mathcal{G} \cdot \left(\frac{\mathcal{N} \times_1 W^T \times_2 H^T \times_3 Q^T}{\mathcal{D} \times_1 W^T \times_2 H^T \times_3 Q^T} \right)^{\gamma(\beta)}, \epsilon \right)$$

/* De même, utilisation de propriétés des produits tensoriels pour réduire la complexité, voir [1] pour le détail. */

• On peut définir de même **Euc-NTD**, par rapport à la norme euclidienne [2] :

$$\arg \min_{W \geq 0, H \geq 0, Q \geq 0, \mathcal{G} \geq 0} \|\mathcal{X} - \mathcal{G} \times_1 W \times_2 H \times_3 Q\|_2^2$$

Pattern musical

Aspect théorique :

Le produit $W \mathcal{G}_{:,b} H^T$ définit des "patterns musicaux" : spectrogrammes à l'échelle de la mesure.

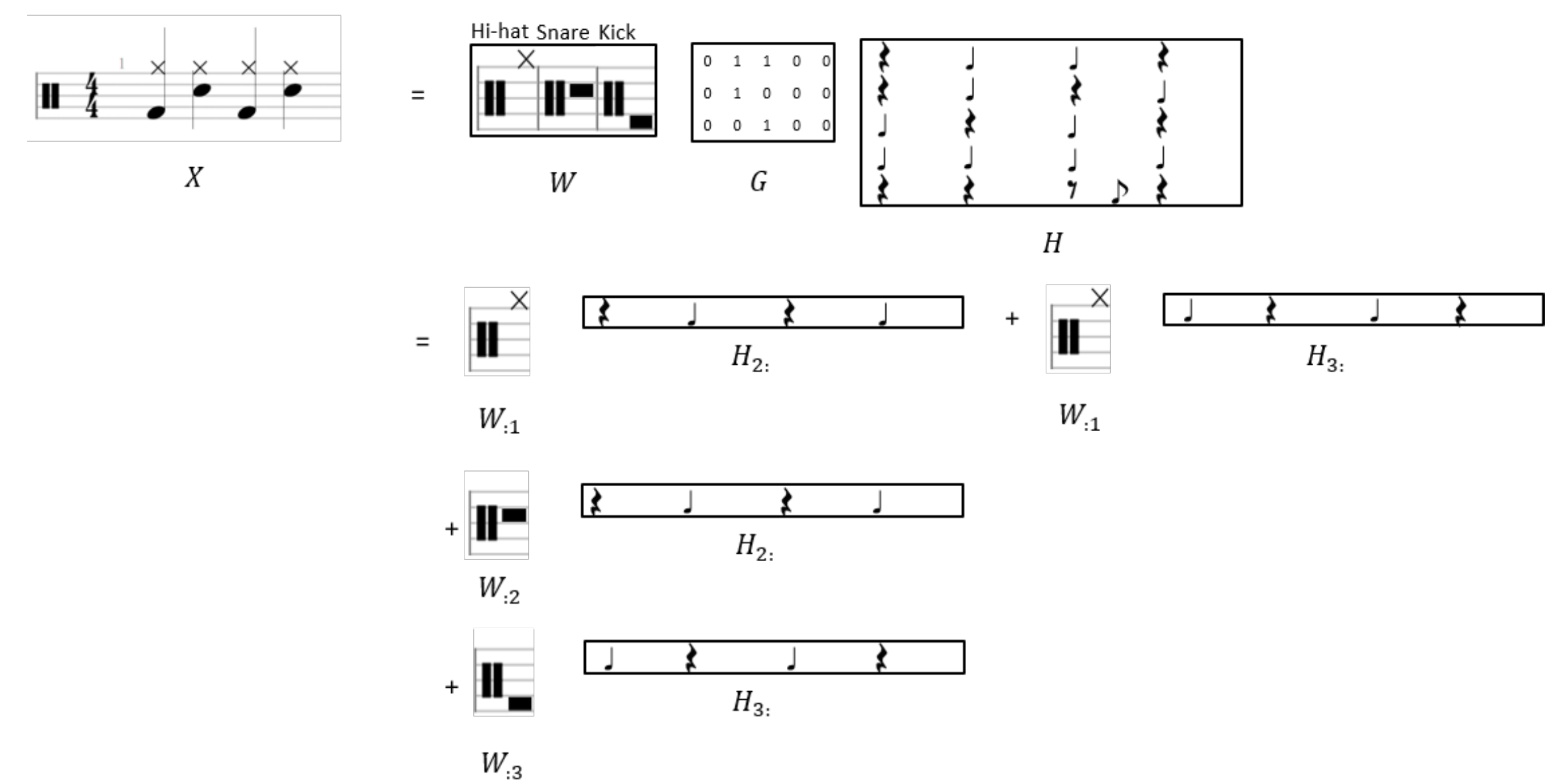


Figure. Un pattern musical théorique (partition de batterie).

Pattern musical en pratique :

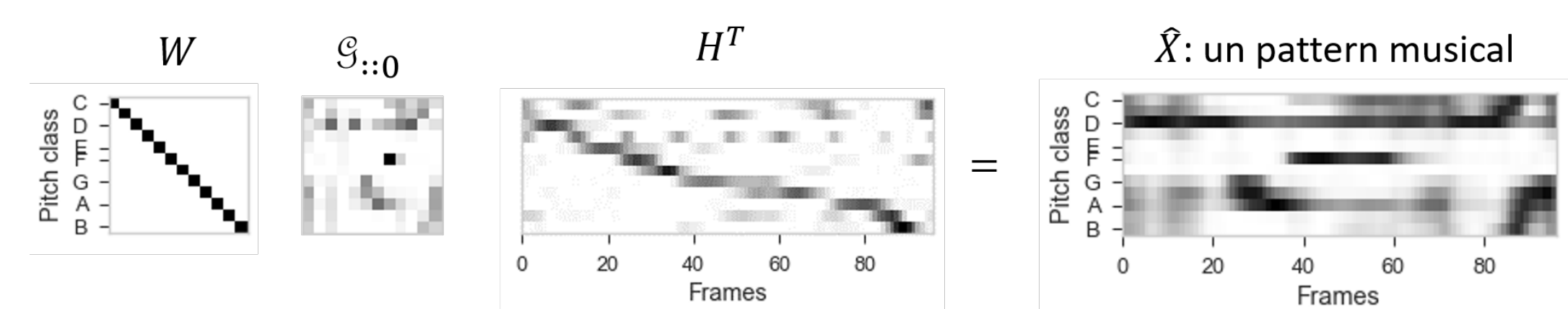


Figure. Un pattern musical, en pratique (chromagramme) : une combinaison linéaire des colonnes de W (extrême-gauche, information fréquentielle) et de H (centre-droit, information rythmique) est définie par la première tranche de \mathcal{G} (centre-gauche), résultant en un chromagramme d'une mesure.

En estimant la phase, on peut écouter les résultats de la décomposition : <https://ax-le.github.io/resources/examples/ListeningNTD.html>



NTD pour la segmentation structurale

Matrice Q^T : patterns musicaux comme descripteurs par mesure.

Chaque spectrogramme, à l'échelle de la mesure, se décompose : $\mathcal{X}_{:,b} \approx \sum_{b'=1}^{B'} Q(b, b') W \mathcal{G}_{:,b'} H^T$

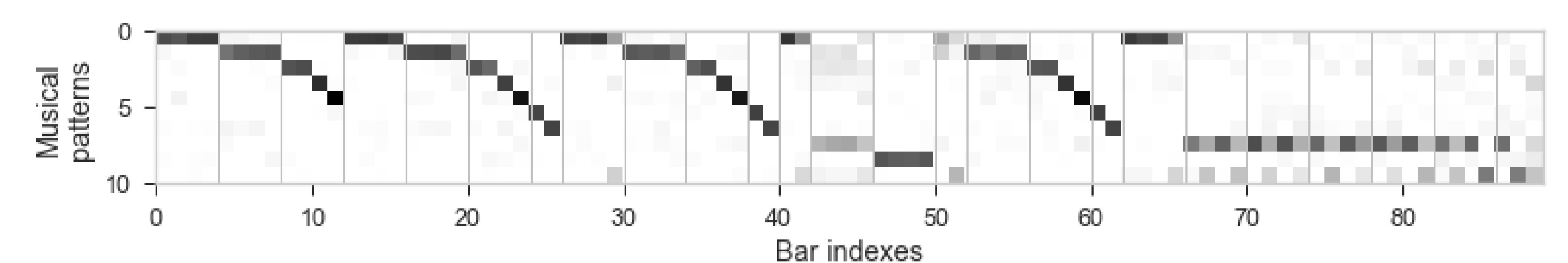


Figure. Matrice Q^T pour "Come Together". Lignes grises : annotations de segmentation.

Performances en segmentation structurale (RWC Pop).

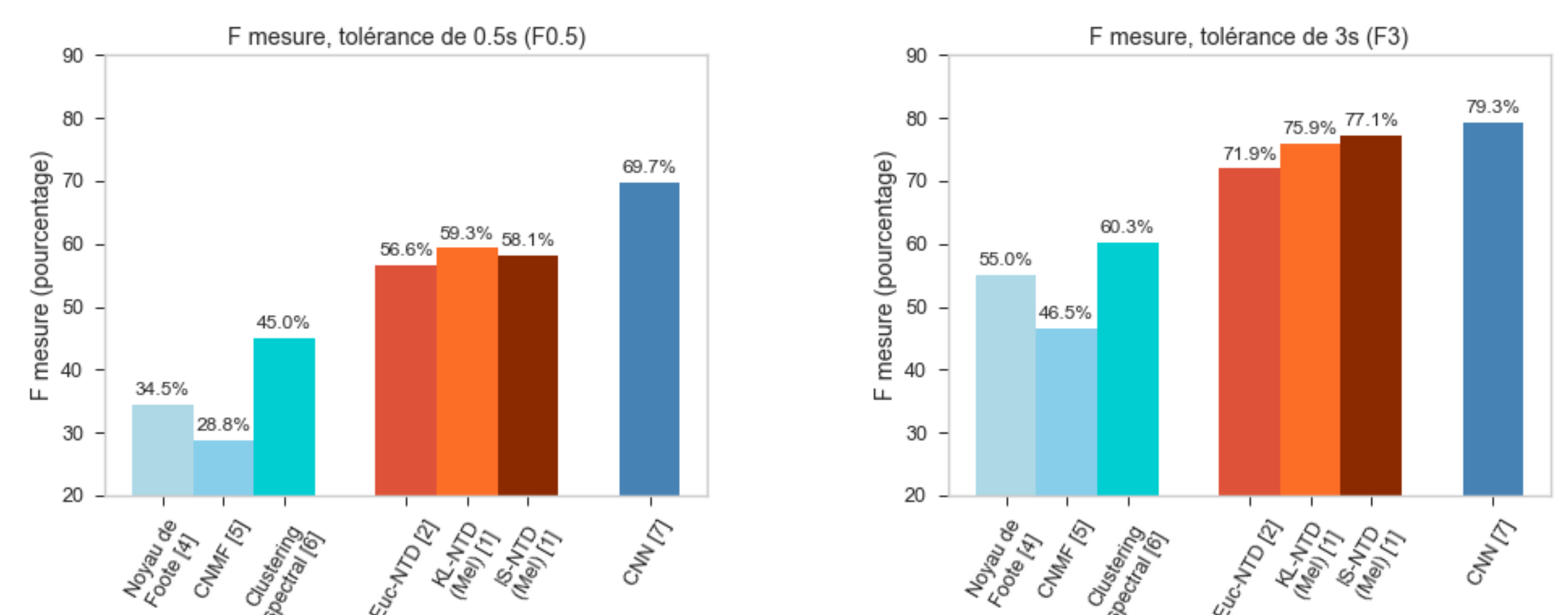


Figure. Comparaison des performances de segmentation (F mesures) avec l'état-de-l'art, respectivement [4, 5, 6] (algorithmes non-supervisés) et [7] (réseau de neurones supervisé).

Articles, code et notebooks

Voir les articles [1, 2]. Le code est en accès libre, avec des tutoriels et notebooks d'expériences :



Code NTD : <https://gitlab.inria.fr/amarmore/nonnegative-factorization>

Code de traitement et expériences : <https://gitlab.inria.fr/amarmore/musicntd>

Références

- A. Marmoret, F. Voorwinden, V. Leplat, J. E. Cohen, and F. Bimbot, "Nonnegative tucker decomposition with beta-divergence for music structure analysis of audio signals," *arXiv preprint arXiv:2110.14434*, 2021.
- A. Marmoret, J. Cohen, N. Bertin, and F. Bimbot, "Uncovering audio patterns in music with nonnegative Tucker decomposition for structural segmentation," in *ISMIR*, pp. 788–794, 2020.
- C. Févotte and J. Idier, "Algorithms for nonnegative matrix factorization with the β -divergence," *Neural computation*, vol. 23, no. 9, pp. 2421–2456, 2011.
- J. Foote, "Automatic audio segmentation using a measure of audio novelty," in *2000 IEEE International Conference on Multimedia and Expo. ICME2000. Proc. Latest Advances in the Fast Changing World of Multimedia (Cat. No. 00TH8532)*, vol. 1, pp. 452–455, IEEE, 2000.
- O. Nieto and T. Jehan, "Convex non-negative matrix factorization for automatic music structure identification," in *2013 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 236–240, IEEE, 2013.
- B. McFee and D. Ellis, "Analyzing song structure with spectral clustering," in *ISMIR*, pp. 405–410, 2014.
- T. Grill and J. Schlüter, "Music boundary detection using neural networks on combined features and two-level annotations," in *ISMIR*, pp. 531–537, 2015.