

INTERFERENCE REDUCTION IN MUSIC RECORDINGS COMBINING KERNEL ADDITIVE MODELLING AND NON-NEGATIVE MATRIX FACTORIZATION

Delia Fano Yela¹, Sebastian Ewert¹, Derry FitzGerald², Mark Sandler¹

Queen Mary University of London, UK¹
Nimbus Centre, Cork Institute of Technology, Ireland²

ABSTRACT

When recording music in a live or studio scenario, unexpected sound events often lead to interferences in the signal. For non-stationary interferences, sound source separation techniques can be used to reduce the interference level in the recording. In this context, we present a novel approach combining the strengths of two sound source separation methods, NMF and KAM. The recent KAM approach applies robust statistics on frames selected by a source-specific kernel to perform source separation. Based on semi-supervised NMF, we extend this approach in two ways. First, we locate the interference in the recording based on detected NMF activity. Second, we improve the kernel-based frame selection by incorporating an NMF-based estimate of the clean music signal. Further, we introduce a temporal context in the kernel, taking musical structure into account. Our experiments show improved separation quality for our proposed method over state-of-the-art methods for interference reduction.

Index Terms— Source separation, Kernel Additive Modelling, Non-Negative Matrix Factorization, Interference Reduction.

1. INTRODUCTION

When recording music professionally, one often has to deal with various types of sound interferences. For example, a person in the audience experiencing a coughing fit during a classical music concert can be a major disturbance. Similarly, fans screaming too close to one of the stage microphones can render the entire channel useless in post-production. Further, studio sessions are often subject to strict time budgets and thus many tracks are only recorded until the sound engineer assesses the last take to be good enough – only to find a door being slammed or an object falling on the floor in this one good take during the actual production.

The difficulty of removing such an interference strongly depends on its type. Stationary interferences, such as mains or fluorescent light hum, can often already be reduced by simple (Wiener) filtering techniques [1]. Non-stationary interferences such as the ones described above, however, require more complex signal models and sound source separation techniques to differentiate noise from non-noise signal components. In this context, Non-Negative Matrix Factorization (NMF) has turned out to be a powerful tool and most state-of-the-art source separation methods are based on NMF variants [2]. The basic idea behind NMF is to model a time-frequency representation of the signal as a product of two matrices. The columns of the first matrix are often interpreted as (spectral) *templates* capturing the spectral properties of the individual sound sources in the signal; the rows of the second matrix are often referred to as the corresponding *activations*, encoding when and how strong each template is active in the input signal.

Applying the original NMF approach [3] to audio and music data, however, was found to rarely yield useful results [4]. Therefore, various extensions were proposed integrating various constraints on the parameter estimation process. Examples include sparsity and temporal continuity constraints [5] or harmonicity constraints [6]. Further, various types of side information have been used, such as user-assisted annotations [7] and musical score information [8]. One of the most widely used and successful approaches is to employ training data (*Supervised NMF*): using recordings containing only a single sound source, corresponding templates representing that source can easily be computed [9]. This way, one can avoid relying on specific statistical independence assumptions for the sources [10]. As a major drawback of this approach, however, the quality of the separation result heavily depends on the assumption that the acoustical conditions in the training material and in the recording to be processed are similar. The more this assumption is violated, the more artefacts are to be expected.

As an alternative to NMF, *Kernel Additive Modelling* (KAM) [11] was proposed for various tasks in source separation, e.g. singing voice separation [12, 13] or the separation of harmonic from percussive signal components [14]. In general, the idea behind KAM is to exploit that the magnitude of a bin in a time-frequency representation is often similar or related to the magnitude of certain other bins – which bins are similar is described by a so called kernel. If the magnitude of a given bin deviates in an unexpected way from the bins defined in the kernel, one can assume that this bin is overlaid by another sound source and we can use the kernel bins to reconstruct the overlaid one. Since some of the kernel bins might be overlaid by other sounds as well, or are not exact repetitions, one uses *Robust Statistics*, in particular order statistics, to identify the commonalities between the bins while neglecting the outliers.

To apply a KAM-based method to a source separation problem, one needs to design a corresponding kernel that identifies similar spectral bins for the sources we want to keep while ignoring the energy associated with other sources. In existing KAM approaches, this kernel design is often rather rudimentary. For example, to eliminate the singing voice from recordings, the methods proposed in [12, 13] assume that the accompaniment playing the harmony changes more slowly than the singing voice and thus that there are many frames with similar accompaniment. The kernel used in [12, 13] is simply a function finding the K most similar frames based on the Euclidean distance. However, using such simple kernels one implicitly assumes that the energy in frames will be dominated by the sound source we want to keep – otherwise the similarity measure fails to identify similar frames. Therefore, while standard KAM is free of the need for suitable training data as in supervised NMF, it might fail to find similar frames if the signal-to-interference ratio is low. In particular, with sudden, loud interferences as to be expected in our application scenario, existing KAM approaches are likely to fail.

Our main idea is to combine the strengths of both methods. In particular, while training data in supervised NMF might not be precise enough to yield a high quality signal model as needed for source separation, it might be discriminative enough to obtain an initial signal model for the music, which can be used to design an adaptive, interference-resilient kernel for KAM. More precisely, we let the user provide keywords to describe the interference (e.g. 'cough') and retrieve corresponding training data from the publicly available freesound¹ archive. After computing templates specific for the interference from the training data, we apply a semi-supervised NMF, i.e. we fix templates for the interference and learn some additional free templates to model the music from the actual input signal. Then, using the (HMM-smoothed) NMF activations for the fixed interference templates, we automatically locate the interference within the recording – this way, in contrast to existing KAM approaches, we can filter the signal only where needed. Second, using the activations for the free templates, we can reconstruct an initial rough estimate for the music, where the interference is strongly reduced as most of the corresponding energy is already captured by the interference templates. Based on this initial model, we identify for each frame affected by the interference a list of similar frames, which are then used within the KAM framework to produce the final output. As additional contributions, we modify the standard kernels used in KAM by incorporating a temporal context into the similarity search which essentially yields a simple regularizer promoting temporal continuity of the kernels across frames, as well as a smoothing technique, which enhances the method's invariance against small variations in the fundamental frequency.

The remainder of the paper is organized as follows. In Section 2 we describe the technical details of our approach. Next, in Section 3 we compare our proposed method with standard KAM and semi-supervised NMF in a series of systematic experiments. Finally, we conclude the paper in Section 4 with an outlook on future work.

2. PROPOSED METHOD

Overall, we develop our method as an extension to *Kernel Additive Modelling (KAM)* [11]. From a modelling point of view, KAM and the more widely known Gaussian Processes (GP) share similar concepts. In both cases, the idea is that for many signals we can estimate the value of a single sample by looking at the value of neighboring samples. For example, a low frequency signal corrupted by white noise can be reconstructed by averaging the values of neighbouring samples. This operation is essentially similar to a low-pass FIR filter, just that KAM and GP enable the use of much more general notions of similarity or neighborhood. KAM differs from GP in several aspects. First, the similarity kernel in KAM can depend on the observations themselves [15], which we exploit in the following. Second, non-Gaussian noise corrupting the sample values can be modelled. Third, as an instance of kernel local regression, KAM does not require the inversion of a data covariance matrix (as in GPs), which typically leads to considerable improvement in terms of computational costs [11].

The KAM framework as a whole is relatively rich, both in possible application scenarios and theory. Due to space constraints, we will only present a smaller subset that was also used in a similar form in the REPET family of methods for singing voice removal [13]. To this end let x be the signal to be processed with $x(t) = s(t) + n(t)$, where s and n are the clean music and the interference signal, respectively. Further, let $X, S \in \mathbb{C}^{F \times T}$ be the spectrograms of x and s . In the following, we exploit that spectral frames in S typically

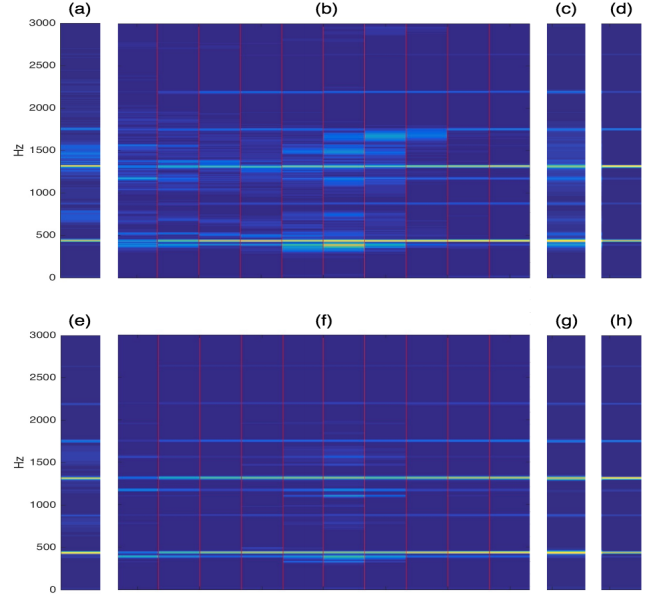


Fig. 1. Individual steps in our proposed method, see text. (a)-(d) using standard KAM, (e)-(h) using our proposed extension.

occur several times in similar form, either because note constellations are repeated over time (as is common in music) or because notes are being held for a while. The interference on the other hand may or may not be repetitive and thus we cannot make any assumptions here. Therefore, we will model only s in KAM without considering the interference n as an actual sound source but just as noise with an unknown distribution. Since s only consists of a single channel, we can eliminate many unnecessary elements in KAM (multi-channel and iterative re-estimation extensions, compare [11]), resulting in a very simple representation. More precisely, let $\mathcal{I} : F \times T \rightarrow F \times T \times K$ be a similarity kernel function that assigns to every time-frequency bin (f, t) a list of K similar bins. As in the case of REPET, we use a frame-wise, K -nearest neighbors (K -NN) function based on the Euclidean distance, i.e. (f, \tilde{t}) is in $\mathcal{I}(f, t)$ if frame \tilde{t} is among the K most similar frames. This process is illustrated in Fig. 1, where for a given frame shown in Fig. 1a the $K = 10$ most similar frames are shown in Fig. 1b.

Once this notion of similarity between bins is established, we can try to calculate a noise-free estimate for each bin (f, t) from the bins in $\mathcal{I}(f, t)$. In KAM this goal is expressed as an optimization problem over so called *model cost functions* \mathcal{L} . More precisely, we get²:

$$\bar{S}(f, t) = \underset{S}{\operatorname{argmin}} \sum_{(f, \tilde{t}) \in \mathcal{I}(f, t)} \mathcal{L}(\bar{X}(f, \tilde{t}), S),$$

where \bar{X} is the magnitude of X . Here \mathcal{L} models our belief regarding how good or bad a specific choice for $\bar{S}(f, t)$ is, considering that we call all elements in $\mathcal{I}(f, t)$ similar to it. A common choice in the KAM framework is $\mathcal{L}(a, b) := |a - b|$. This choice is interesting for two reasons. First, it expresses from a probabilistic point of view that we expect some larger deviations in the difference a and

²Note that the formal requirement to have noisy data extracts for S (as used in KAM) are directly given here as $X(f, \tilde{t})$. This is the result of having only a single sound source in a single channel, which makes the iterative re-estimation in KAM unnecessary and eliminates all elements related to the estimation of the mixing matrix, compare [11].

¹<https://www.freesound.org/>

b , and that this distance is not Gaussian distributed (otherwise a Euclidean distance would be optimal here). Second, this choice leads us to the use of robust statistics in the form of the median, which as an operator is invariant against outliers and thus allows robust parameter estimation in the presence of noise. More precisely, with $\mathcal{L}(a, b) := |a - b|$ the solution to the above problem is:

$$\bar{S}(f, t) = \text{median}(\bar{X}(f, \tilde{t}) | (f, \tilde{t}) \in \mathcal{I}(f, t)).$$

The result is shown in Fig. 1c.

Comparing the results of this approach shown in Fig. 1c with the clean signal in Fig. 1d, we can observe an example of when this approach fails. In particular, comparing Fig. 1a and Fig. 1d, we see that the input frame is overlaid by a strong interference. With such a low signal-to-noise ratio (SNR), the interference dominates in the similarity search based on the Euclidean distance and the kernel function $\mathcal{I}(f, t)$ points to too many noisy examples, which even the median operation cannot eliminate. In particular, despite being strongly invariant against outliers from a robust statistics point of view, the outliers cannot be identified anymore based on a selection of frames as shown in Fig. 1c.

Our idea is now to improve the K -NN search in KAM in several ways, making the kernel function more invariant against the interference signal. To this end, we build a first initial signal model based on NMF using training data. While the training data might differ from the actual interference signal, and thus an actual source separation based on this method would yield results of low quality, it might be good enough to gather more information about the signal and reduce the influence of the interference. More precisely, similar to [16] we let the user provide keywords to describe the interference (e.g. 'cough') and retrieve corresponding example recordings from the freesound archive. Concatenating these recordings into a single file, we compute its magnitude spectrogram \bar{X}_I as well as an NMF factorization $\bar{X}_I \approx W_I H_I$ using the well-known Lee-Seung NMF updates for the generalized Kullback-Leibler divergence D_{KL} [3], i.e. we minimize $D_{\text{KL}}(\bar{X}_I, W_I H_I)$ over non-negative matrices W_I and H_I . The only parameter here is the NMF rank R_1 . After this, the columns of W_I contains templates reflecting the spectral properties of the interference signal.

In a next step, we employ NMF to model our input spectrogram \bar{X} using combinations of interference templates and music templates. Here, the interference templates can be kept fixed and we only need to learn the music templates, which is often referred to as *semi-supervised NMF*. More precisely, we minimize the function $D_{\text{KL}}(\bar{X}, W_I H_I + W_M H_M)$ over H_I , W_M and H_M (i.e. we fix W_I). In this case, the update rules are similar to regular NMF:

$$H_I \leftarrow H_I \odot \frac{W_I^\top \mathcal{R}}{W_I^\top \cdot J} \quad \text{and} \quad H_M \leftarrow H_M \odot \frac{W_M^\top \mathcal{R}}{W_M^\top \cdot J},$$

$$W_M \leftarrow W_M \odot \frac{\mathcal{R} H_M^\top}{J \cdot H_M^\top}, \quad \text{with} \quad \mathcal{R} := \frac{\bar{X}}{W_I H_I + W_M H_M}$$

and J the all-one matrix. After convergence, the rows of H_I capture the activations of the interference templates, while $W_M H_M$ yields an approximation of the magnitude spectrogram for the music. Using these two interpretations, we employ these results for two different purposes. First, we use H_I to identify where the interference is, which will enable us to filter only frames with interference (in contrast to regular KAM). To this end, we sum the values in H_I in each frame to obtain a single curve indicating interference activity, which we decode using an HMM, resulting in a binary, frame-wise interference indicator vector I . The parameters of the HMM implemented, detection threshold and cost of changing state, were adjusted to favour recall over precision in the detection.

Next, we exploit that while the interference templates might not perfectly reflect the properties of the target interference (and thus a separation based on this model would be of low quality), they do capture typically a considerable amount of interference energy in the signal. Therefore, we can improve the K -NN search in KAM kernel if we replace the input spectrogram \bar{X} containing the interference with the NMF approximation for the music $\tilde{X} := W_M H_M$. The resulting improvement is clearly visible in Fig. 1. If we replace the \bar{X} -frame (Fig. 1a) with the corresponding \tilde{X} -frame (Fig. 1e), we see that the frames selected as nearest neighbours (Fig. 1f) are much closer to the actual target (Fig. 1d = Fig. 1h). The median filter can then remove remaining noise robustly, bringing the result (Fig. 1g) much closer to the target (Fig. 1h).

However, in the case of little or no musical pattern repetition in the mixture, the frames selected as nearest neighbours in \tilde{X} could contain a significant amount of interference energy in \bar{X} and would, therefore, render the median filtering ineffective as the interference would no longer be an outlier. In order to address this scenario, we propose to check if the frames selected as nearest neighbours were also previously identified as interference frames and, in that case, use the \tilde{X} -frames for median filtering instead of the \bar{X} -frames. As it is shown in Section 3, this contribution will reduce the interference impact on the median filtering output.

Despite these improvements in K -NN search, a second problem we observed is that the kernel \mathcal{I} was often changing considerably between frames in the sense that often $(f, \tilde{t}) \in \mathcal{I}(f, t)$ would not imply $(f, \tilde{t} + 1) \in \mathcal{I}(f, t + 1)$. Without this property, however, we observed a slight pitch jitter in the magnitude across frames after median filtering, which did not fit in the signal musical pattern and was audible in the resulting time domain. To further temporally stabilize the kernel function, we propose incorporating a temporal context into the similarity search. More precisely, instead of comparing frames t and \tilde{t} with a simple Euclidean distance $\sum_f (\tilde{X}(f, t) - \tilde{X}(f, \tilde{t}))^2$, we employ

$$\sum_f \sum_{c=-C}^C (\tilde{X}(f, t+c) - \tilde{X}(f, \tilde{t}+c))^2$$

as frame distance in the K -NN search, where C specifies the *temporal extent*. We found this simple extension to act as a surprisingly effective temporal regularizer for \mathcal{I} . Further, we found that filtering \tilde{X} slightly in frequency direction before the K -NN search using a small Gaussian kernel additionally improved the results, as it makes the similarity search invariant to small changes in the fundamental frequency of harmonic sounds.

To perform the actual separation, we employ soft masking (Wiener filtering). In particular, our method yields an estimate \bar{S} for the magnitude spectrogram of the music. We define a corresponding estimate for the noise, here interference, as $\bar{N} = \max(\bar{X} - \bar{S}, 0)$. This way, we can obtain an estimate S for the complex music spectrogram via $S = \frac{\bar{S}}{\bar{N} + \bar{S}} \odot X$. Overall, we found our method combining NMF and KAM to improve over both approaches considerably, which we demonstrate in the next section.

3. EXPERIMENTS

We evaluated our proposed method using freely available recordings, in particular interferences and instrumental solo stems from multi-track recordings [17]. We chose interferences that typically occur in a live or studio scenario including cough sounds, door slams, sounds of objects of different material being dropped, chair-drag sounds as well as audience screams. The music corresponds to 58 instrumental

	NSDR			NSIR		
	0dB	-3dB	-6dB	0dB	-3dB	-6dB
Prop.	6.95	5.98	3.61	15.06	15.68	14.10
NMF	4.98	3.40	1.30	13.07	14.26	15.23

Table 1. Comparison of our methodology with the state-of-the-art for different SNR values

mono stems from the multitrack MedleyDB dataset [17], covering 23 different instruments ranging from guitar, violin, piano over to bass, trombone or flute.

Recordings of interferences were sourced freesound.org to potentially render the method more accessible for the public as described in [16]. However, this implies that the quality and number of training samples is subject to the contributions available on the net, and therefore explains why, in our case, each interference has a different amount of training data, ranging from 10 clean scream samples to 40 clean coughs tracks. The separation quality is expected to improve as the number of tracks in the training data increases.

We created test recordings by making artificial linear mixes of stems and test interference recordings independent from the training data and of each other (other acoustic conditions). In order to achieve a controlled mix of instrumental and interference levels, all of the tracks used here are normalised to a certain RMS energy. Then three interferences are added to the music at different SNR, measured on the segment where the interference is active. The final mix is a 30s long mono instrumental track with three sounds of the same kind interfering at different times at a certain SNR.

We therefore tested the proposed method for the 290 mixtures of the 58 instrumental stems with the 5 types of interferences, and we evaluated the separation performance with BSS Eval toolbox [18], obtaining a Signal-to-Distortion Ratio (SDR) and Signal-to-Interference Ratio (SIR) for each mixture separation. The SDR is used as a measure to indicate the overall separation performance, whereas the SIR shows how much of the interference signal is left in the signal estimate. To indicate the improvement over the raw music-interference mix, we employ the normalized SDR/SIR (NSDR/SIR) as in [19], i.e. from the SDR for our method we subtract the SDR for the mix. This way, we can account for the fact that a separation at a low SNR is more difficult than at a high SNR, making results for different SNRs more comparable.

Here we have chosen supervised-NMF to represent the current state-of-the-art method to quantitatively compare its separation performance to the proposed technique’s one. In order to obtain a competitive baseline, we use the same learned dictionary for both methods and we also optimise the NMF rank value with a parameter sweep. In Table 1 and Table 2 we can find the overall results for the separation performance of the proposed method, averaged across all NSDR/ NSIR values of every mixture, as well as the results for the semi-supervised NMF approach.

In view of the results presented on Table 1, the proposed method achieves better separation performance than the state-of-the-art method not only for a 0dB SNR mixture, but also for mixtures where the interference is 3dB and 6dB below the instrumental RMS energy. Although the overall performance drops with the interference energy contribution, our method improves the state-of-the-art, here over NMF, especially when the separation is more difficult and more likely to be required.

In order to measure the influence of the individual components of the proposed method, Table 2 shows the results separately for every

	NSDR	NSIR
Variant 1	7.67	14.60
Variant 2	8.54	16.67
Variant 3	9.58	15.59

Table 2. Contributions impact on separation results for 0dB SNR

contribution presented in this paper, in form of variations of the main method. To provide another angle on the results and focus on the positions where the interferences actually happen, we evaluated the separation performance by averaging across the three segments in the mix where the interference is active, and so the resultant NSDR scores are higher than those obtained on the overall signals.

Starting with a baseline KAM methodology, analogous to the one described in [12], *Variant 1* adds the NMF interference detection step introduced earlier in this paper. The high NSDR shows the interference was successfully identified and reduced.

Variant 2 algorithm represents the improved similarity measure of our proposed method. The frame-wise similarity search of *Variant 1* now also includes a temporal context in *Variant 2* with the additional group-wise similarity measure performed on the interference-cleaned NMF instrumental estimate instead on the original mixture, as in *Variant 1*. The higher NSDR shows that the temporal context in the similarity measure not only stabilizes the kernel and results perceptually, but also quantitatively. Nonetheless it is important to point out the limited pattern repetition scenario due to the use of short music snippets of only 30 seconds. It is expected to see an NSDR improvement when the mixtures contain a whole song, and therefore, more pattern repetition. In addition, we can now also objectively state that being accurate in measuring and defining close frames has a major impact on the overall separation performance and source identification.

Variant 3 is an extension of *Variant 2* incorporating an additional smoothing filter across the close frames that are now extracted, not exclusively from the original signal, but also from the interference-cleaned NMF instrumental estimate for those close frames who were previously identified as interference-active. *Variant 3* improves the NSDR as there is no longer a loss of energy due to the estimate’s pitch misalignment in *Variant 2*. However, the NSIR values can be lower in this case where the smoothing pitch invariance filter may have deteriorated the precision of the source spectral characterisation.

Overall, these results show that our proposed extensions improve the interference and instrument source separation in an objective manner.

4. CONCLUSION

We have presented a new method for interference reduction combining NMF and KAM. Our method exploits the advantages of both techniques, using a spectral dictionary to detect the interferences occurrences as well as to produce an initial instrumental estimate using NMF. This estimate is used to improve the similarity measure of KAM, taking also a temporal musical context into account. The methodology presented shows an improvement of the state-of-the-art for solo instrumental signals. In addition, the processing only takes place for those frames affected by the interference, yielding more pleasant audio quality results that could potentially be used in a final mix. Possible avenues for extending this work would include an improved similarity search as well as the implementation of source-specific kernels both in time and frequency direction.

5. REFERENCES

- [1] Monson H. Hayes, *Statistical Digital Signal Processing and Modeling*, Wiley, 1st edition, 1996.
- [2] Andrzej Cichocki, Rafal Zdunek, and Anh Huy Phan, *Non-negative Matrix and Tensor Factorizations: Applications to Exploratory Multi-Way Data Analysis and Blind Source Separation*, John Wiley and Sons, 2009.
- [3] Daniel D. Lee and H. Sebastian Seung, “Algorithms for non-negative matrix factorization,” in *Proceedings of the Neural Information Processing Systems (NIPS)*, Denver, CO, USA, 2000, pp. 556–562.
- [4] Derry FitzGerald, Matt Cranitch, and Eugene Coyle, “Extended nonnegative tensor factorisation models for musical sound source separation (article id 872425),” *Computational Intelligence and Neuroscience*, vol. 2008, 2008.
- [5] Tuomas Virtanen, “Monaural sound source separation by nonnegative matrix factorization with temporal continuity and sparseness criteria,” *IEEE Transactions on Audio, Speech and Language Processing*, vol. 15, no. 3, pp. 1066–1074, 2007.
- [6] Nancy Bertin, Roland Badeau, and Emmanuel Vincent, “Enforcing harmonicity and smoothness in Bayesian non-negative matrix factorization applied to polyphonic music transcription,” *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 18, no. 3, pp. 538–549, 2010.
- [7] Paris Smaragdis, “User guided audio selection from complex sound mixtures,” in *Proceedings of the ACM Symposium on User Interface Software and Technology (UIST)*, New York, NY, USA, 2009, pp. 89–92.
- [8] Sebastian Ewert, Bryan Pardo, Meinard Müller, and Mark D. Plumbley, “Score-informed source separation for musical audio recordings: An overview,” *IEEE Signal Processing Magazine*, vol. 31, no. 3, pp. 116–124, May 2014.
- [9] Yong-Choon Cho and Seungjin Choi, “Learning nonnegative features of spectro-temporal sounds for classification,” in *Proceedings of InterSpeech*, 2004.
- [10] Samer Abdallah and Mark Plumbley, “Polyphonic transcription by non-negative sparse coding of power spectra,” in *Proceedings of the International Society for Music Information Retrieval Conference (ISMIR)*, Barcelona, Spain, 2004, pp. 318–325.
- [11] Antoine Liutkus, Derry FitzGerald, Zafar Rafii, Bryan Pardo, and Laurent Daudet, “Kernel additive models for source separation,” *IEEE Transactions on Signal Processing*, vol. 62, no. 16, pp. 4298–4310, 2014.
- [12] Derry FitzGerald, “Vocal separation using nearest neighbours and median filtering,” in *Proceedings of the Irish Signals and Systems Conference (ISSC)*, 2012, pp. 1–5.
- [13] Zafar Rafii and Bryan Pardo, “Repeating pattern extraction technique (REPET): A simple method for music/voice separation,” *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 21, no. 1, pp. 71–82, 2013.
- [14] Derry FitzGerald, “Harmonic/percussive separation using median filtering,” in *Proceedings of the International Conference on Digital Audio Effects (DAFx)*, Graz, Austria, 2010, pp. 246–253.
- [15] Joaquin Quiñonero-Candela and Carl Edward Rasmussen, “A unifying view of sparse approximate gaussian process regression,” *Journal of Machine Learning Research*, vol. 6, no. Dec, pp. 1939–1959, 2005.
- [16] Dalia El Badawy, Ngoc QK Duong, and Alexey Ozerov, “On-the-fly audio source separation,” in *Proceedings of the IEEE International Workshop on Machine Learning for Signal Processing (MLSP)*, 2014, pp. 1–6.
- [17] Rachel M Bittner, Justin Salamon, Mike Tierney, Matthias Mauch, Chris Cannam, and Juan Pablo Bello, “Medleydb: A multitrack dataset for annotation-intensive mir research,” in *ISMIR*, 2014, pp. 155–160.
- [18] Cédric Févotte, Rémi Gribonval, and Emmanuel Vincent, “Bss_eval toolbox user guide–revision 2.0,” 2005.
- [19] Antoine Liutkus, Derry Fitzgerald, and Zafar Rafii, “Scalable audio separation with light kernel additive modelling,” in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2015, pp. 76–80.