

# SOURCE SEPARATION WITH SCATTERING NON-NEGATIVE MATRIX FACTORIZATION

*Joan Bruna, Pablo Sprechmann, Yann LeCun*

New York University  
Courant Institute of Mathematical Sciences  
{bruna, pablo}@cims.nyu.edu, yann@cs.nyu.edu

## ABSTRACT

This paper presents a single-channel source separation method that extends the ideas of Nonnegative Matrix Factorization (NMF). We interpret the approach of audio demixing via NMF as a cascade of a pooled analysis operator, given for example by the magnitude spectrogram, and a synthesis operators given by the matrix decomposition. Instead of imposing the temporal consistency of the decomposition through sophisticated structured penalties in the synthesis stage, we propose to change the analysis operator for a deep scattering representation. This new signal representation is invariant to smooth changes in the signal, consistent with its temporal dynamics. We evaluate the proposed approach in a speech separation task obtaining promising results.

**Index Terms**— source separation, scattering, non-negative matrix factorization.

## 1. INTRODUCTION

The problem of source separation has been widely studied in the speech processing community [1, 2]. It becomes particularly challenging when only one microphone is used, or in the presence of non-stationary background noise, which is a very common situation in many applications encountered, e.g., in telephony. We approach this problem as a monaural source separation method by modeling the speech at an appropriate temporal resolution.

The decomposition of time-frequency representations, such as the power or magnitude spectrogram in terms of elementary atoms of a dictionary, has become a popular tool in audio processing. Non-negative matrix factorization (NMF) [3], have been widely adopted in various audio processing tasks, including in particular source separation, see [4] for a recent review. There are many works that follow this line in speech separation [5, 6] and enhancement [7, 8].

In plain NMF, signals that can be well approximated with the learned dictionary are likely to resemble the training data on a frame by frame manner. They might, however, not be temporally consistent at larger temporal scales. Standard NMF approaches treat different time-frames independently, ignoring the temporal dynamics of the signals. In order to

overcome this limitation, many works have proposed regularized extensions of NMF to promote learned structure in the codes. Examples of these approaches are, temporal smoothness of the activation coefficients [9], including co-occurrence statistics of the basis functions [10], and learned temporal dynamics [11, 12, 13].

NMF-based source separation methods can be thought as the concatenation of two operators. First, the signal is represented in a feature space given by a non-linear analysis operator, typically defined as the magnitude of a time-frequency representation such as the Short-Time Fourier Transform (STFT). Then a synthesis operator, given by the dictionary learning stage, is applied to produce an unmixed in the feature space. The separation is obtained by inverting these representations. Performing the separation in the non-linear representation is key to the success of the algorithm. The magnitude STFT is in general sparse (simplifying the separation process) and invariant to variations in the phase, thus relieving the NMF from learning this irrelevant variability. This comes at the expense of inverting the unmixed estimates in the feature space, normally known as the phase recovery problem [14]. In the case of standard NMF, this is typically done via Wiener filtering.

In this work, rather than optimizing a coding scheme with improved temporal coherence, we concentrate in the extraction of discriminative and stable features. For that purpose, it is crucial to increase the temporal context of the representation, reducing uninformative variability while preserving distinctive speech characteristics. Increasing the temporal scale of STFT or MEL representations results in loss of important discriminative information [15]. In order to overcome the limitations of these shallow representations, scattering transforms [15, 16] cascade several stages of complex wavelet decompositions and complex modulus, yielding discriminative representations with the ability to capture temporal structures at larger scales, e.g. smooth changes in pitch and envelope. Scattering transforms achieve state-of-the-art results on auditory texture discrimination, and music genre recognition [15, 17]. Recently they have been considered in the setting of blind source separation in [18], but here we concentrate in the supervised (and semi-supervised) framework. A dictionary learnt to represent the signal in this deep repre-

sensation implicitly learns the short term temporal dynamics of the signal, capitalizing on the stability properties of scattering coefficients [16].

Our claim is that an important part of the consistency that is imposed via structured NMF, can be eliminated with a better signal representation. In this new setting one can learn the temporal dynamics with a very simple NMF encoding. However, the problem that becomes more difficult is that of inverting the non-linear representation. Recent studies in textured sound synthesis from scattering moments have solved this problem successfully using gradient descent algorithms [19]. Synthesis models with coherent dictionaries are known to be highly unstable representations [20]. Thus, training them to satisfy slowness and temporal consistency can be challenging. In contrast, analysis operators are stable by construction.

Section 2 reviews non-negative matrix factorization, while Section 3 describes scattering representations for speech. Our source separation algorithm is described in Section 4 and numerical experiments on TIMIT and GRID datasets are reported in Section 5.

## 2. NMF SPEECH SOURCE SEPARATION

We consider the setting in which we are given a temporal signal  $x(t)$  that is the sum of two speech signals  $x_i(t)$ ,  $i = 1, 2$ :

$$x(t) = x_1(t) + x_2(t), \quad (1)$$

and we aim at finding estimates  $\hat{x}_i(t)$ . NMF-based source separation techniques typically operate on a non-negative time-frequency representation of  $x(t)$ , such as the spectrogram or the power spectrum, that we denote as  $\Phi(x) \in \mathbb{R}^{m \times n}$ , comprising  $m$  frequency bins and  $n$  temporal frames. NMF attempts to find the non-negative activations  $Z_i \in \mathbb{R}^{q \times n}$ ,  $i = 1, 2$  best representing the different speech components in two dictionaries  $D_i \in \mathbb{R}^{m \times q}$ . This task is achieved through the solution of

$$\min_{Z_i \geq 0} \mu(\Phi(x) | \sum_{i=1,2} D_i Z_i) + \lambda \sum_{i=1,2} \mathcal{R}(Z_i). \quad (2)$$

The first term in the optimization objective measures the dissimilarity between the input data and the estimated channels. Frequent choices of  $\mu$  are the squared Euclidean distance, the Kullback-Leibler divergence, and the Itakura-Saito divergence, for which there exist standard optimization algorithms [21]. In this work we concentrate on a reweighted Euclidean distance, but any other option could be used instead. The second term in the minimization objective is included to promote some desired structure of the activations. Once the optimal activations are solved for, the spectral envelopes of the speech and the noise are estimated as  $D_i Z_i$ . Since these estimated speech spectrum envelope contain no phase information, they are used to build soft masks to filter the mixture signal [22].

## 3. SCATTERING TRANSFORM

Discriminative features having longer temporal context can be constructed with the scattering transform [15, 16]. This section reviews its definition and main properties when applied to speech signals.

### 3.1. Wavelet Filter Bank

A wavelet  $\psi(t)$  is a band-pass filter with good frequency and spatial localization. We consider a complex wavelet with a quadrature phase, whose Fourier transform satisfies  $\mathcal{F}\psi(\omega) \approx 0$  for  $\omega < 0$ . We assume that the center frequency of  $\mathcal{F}\psi$  is 1 and that its bandwidth is of the order of  $Q^{-1}$ . Wavelet filters centered at the frequencies  $\lambda = 2^{j/Q}$  are computed by dilating  $\psi$ :  $\psi_\lambda(t) = \lambda \psi(\lambda t)$ , and hence  $\mathcal{F}\psi_\lambda(\omega) = \hat{\psi}(\lambda^{-1}\omega)$ . We denote by  $\Lambda$  the index set of  $\lambda = 2^{j/Q}$  over the signal frequency support, with  $j \leq J$ , and we impose that these filters fully cover the positive frequencies:

$$\forall \omega \geq 0, \quad 1 - \epsilon \leq |\mathcal{F}\phi(\omega)|^2 + \frac{1}{2} \sum_{\lambda \in \Lambda} |\mathcal{F}\psi_\lambda(\omega)|^2 \leq 1, \quad (3)$$

for some  $\epsilon < 1$ , where  $\phi(t)$  is the lowpass filter carrying the low frequency information at scales larger than  $2^J$ . The resulting filter bank has a constant number  $Q$  of bands per octave. The wavelet transform of a signal  $x(t)$  is

$$Wx = \{x * \phi(t), x * \psi_\lambda(t)\}_{\lambda \in \Lambda}.$$

Thanks to (3), one can verify that

$$\|x\|^2(1 - \epsilon) \leq \|x * \phi\|^2 + \sum_{\lambda \in \Lambda} \|x * \psi_\lambda\|^2 \leq \|x\|^2. \quad (4)$$

### 3.2. Joint Time-Frequency Scattering

Scattering coefficients provide a nonlinear representation computed by iterating over wavelet transforms and complex modulus nonlinearities. First order scattering coefficients are local averages of wavelet amplitudes:

$$\forall \lambda \in \Lambda, \quad Sx(t, \lambda) = |x * \psi_\lambda| * \phi(t).$$

The Q-factor  $Q_1$  adjusts the frequency resolution of these wavelets, and for speech it is typically around  $Q_1 = 32$ . Due to the temporal average, first order scattering coefficients provide no information on the time variation of the scalogram  $|x * \psi_{\lambda_1}(t)|$  at temporal scales smaller than  $2^J$ . It averages all modulations and transient events, and thus loses perceptually important information.

Second order scattering coefficients recover information on audio modulations and pitch temporal variations by computing the wavelet coefficients of the envelopes  $|x * \psi_{\lambda_1}|$ , and their local averages:

$$\forall \lambda_2, \quad Sx(t, \lambda_1, \lambda_2) = \||x * \psi_{\lambda_1}| * \psi_{\lambda_2}\| * \phi(t).$$

These multiscale variations of each envelope  $|x * \psi_{\lambda_1}|$  specify the amplitude modulations of  $x(t)$  [15] and thus have the capacity to detect rhythmic structures appearing at different frequency bands. The Q-factor  $Q_2$  of the second family of wavelets  $\psi_{\lambda_2}$  controls the time-frequency resolution of the transform. Since the envelopes  $|x * \psi_{\lambda}|$  have bandwidth  $\sim 2^{-j}Q_1^{-1}$ , one typically chooses dyadic  $Q_2 = 1$  second order wavelets. Scattering coefficients have a negligible amplitude for  $\lambda_2 > \lambda_1$  because  $|x * \psi_{\lambda_1}|$  is a regular envelop whose frequency support is below  $\lambda_2$  [16]. Scattering coefficients are thus computed only for  $\lambda_2 < \lambda_1$ .

Scattering transforms have been extended along the frequency variables to capture the joint time-frequency variability of spectral envelopes and therefore provide representations locally stable to pitch variations [15]. We denote  $\gamma = \log_2 \lambda_1$ , and consider the scalogram as a two-dimensional function of  $\gamma$  and  $t$ :

$$F(\gamma, t) = |x * \psi_{2^\gamma}(t)|.$$

In this work, we consider a second layer scattering with a separable wavelet transform  $F * \Psi_{\gamma_2, \lambda_2}(\gamma, t)$ , with

$$\Psi_{\gamma_2, \lambda_2}(\gamma, t) = \tilde{\psi}_{\gamma_2}(\gamma) \psi_{\lambda_2}(t).$$

The temporal wavelets  $\psi(t)$  are dyadic complex Morlet wavelets. In this implementation, we choose  $\tilde{\psi}$  to be dyadic real Haar wavelets to preserve good frequency localization.

The resulting second order scattering coefficients are thus

$$\tilde{S}x(t, \lambda_1, \gamma_2, \lambda_2) = |F * \bar{\Psi}_{\gamma_2, \lambda_2}| * \bar{\phi}(\log_2 \lambda_1, t),$$

where  $\bar{\phi}(\gamma, t)$  is a two-dimensional blurring kernel with temporal scale  $2^J$  and log-frequency scale  $2^{J_h}$ . The final representation regrouping first and second order scattering coefficients and sampling at intervals  $k 2^{J-\delta}$  is  $\Phi(x) = \{Sx(k 2^{J-\delta} k, \lambda_1), \tilde{S}x(k 2^{J-\delta}, \lambda_1, \gamma_2, \lambda_2)\}$ , where the oversampling factor  $\delta$  is typically set to  $\delta = 1$ .

#### 4. SOURCE SEPARATION ALGORITHM

We show in this section how the inverse problem of source separation can be solved via a sparse NMF in the scattering domain, followed by phase recovery. We consider the supervised monoaural source separation problem (1), in which the components  $x_i$ ,  $i = 1, 2$  come from sources for which we have training data  $X_i = \{x_{ij}\}_{j \leq K}$ , and one is asked to produce estimates  $\hat{x}_i$ .

The supervision provides prior information on the nature of each of the components. However, high-dimensional speech signals have large variability, most of which is uninformative for the purposes of estimating  $x_i$  in (1). The training data can be exploited more efficiently in the scattering domain, since intra-class variability given by small pitch and timber variations is linearized up to temporal scales  $2^J$  without loosing as much discriminative information as the spectrogram [15, 16].

Let  $\Phi(X_i)$  be the scattering representation of the training examples of each source. We consider a non-linear approximation of each source using sparse NMF:

$$\min_{D_i \geq 0, Z_i \geq 0} \frac{1}{2} \|\Phi(X_i) - D_i Z_i\|_F^2 + \lambda \|Z_i\|_1. \quad (5)$$

This model exploits the linearization properties of scattering coefficients since it searches low-dimensional linear approximations.

At test time, given an input  $x$ ,  $x_1$  and  $x_2$  are estimated by minimizing

$$\min_{x'_i, z_i} \sum_{i=1,2} \frac{1}{2} \|\Phi(x'_i) - D_i z_i\|_2^2 + \lambda \|z_i\|_1 \quad s.t. \quad x = x'_1 + x'_2. \quad (6)$$

Problem (6) is minimized with an alternating gradient descent between  $x'_i$  and  $z_i$ . Fixing  $z_i$  and minimizing with respect to  $x'_i$  requires locally inverting the scattering operator  $\Phi$ , which amounts to solve an overcomplete phase recovery problem and can be solved with gradient descent, as shown in [19]. Fixed  $x'_i$ , solving for  $z_i$  is a standard  $\ell_1$  non-negative sparse coding problem, which can be solved efficiently with proximal splitting algorithms. In this work, we use the LARS algorithm using the SPAMS package [23].

When the analysis operator  $\Phi$  is able to produce sparse representations of the sources, then

$$\begin{aligned} \|\Phi(x'_1) - D_1 z_1\|_2^2 + \|\Phi(x'_2) - D_2 z_2\|_2^2 &\approx \\ \|\Phi(x'_1) + \Phi(x'_2) - \sum_{i=1,2} D_i z_i\|_2^2 &\approx \|\Phi(x) - D_1 z_1 - D_2 z_2\|_2^2, \end{aligned}$$

which can be used in practice to produce a greedy initialization for (6) as follows. We first obtain  $\widehat{\Phi(x_i)} = D_i z_i^*$ , where the  $z_i^*$  are defined as

$$z_i^* = \arg \min_{z_i} \frac{1}{2} \|\Phi(x) - \sum_{i=1,2} D_i z_i\|_2^2 + \lambda \|z_i\|_1.$$

Since the scattering satisfies  $\Phi(x) = \{A|W_1 x|, A|W_2|W_1 x|\}$ , where  $A$  is the lowpass filter and  $W_1$  and  $W_2$  are respectively the first and second layer wavelet decompositions, we can produce an estimate  $\hat{x}_i$  from  $\widehat{\Phi(x_i)}$  by using the complex phases of  $W_1 x$  and  $W_2|W_1 x|$ .

#### 5. EXPERIMENTAL RESULTS

**Evaluation settings.** We evaluated the proposed method in two settings: speaker-specific and multi-speaker. In the first setting we trained a speaker-specific model for each speaker in the mixture and tested it using sentences (from the same speakers) outside the training set. In the second setting, we trained a generic model on a mixed group of male and female speakers, none of which were included in the test set. All signals were mixed at 0 dB and clips resampled to 16 KHz.

	Speaker-Specific			Multi-Speaker		
	SDR	SIR	SAR	SDR	SIR	SAR
NMF	8.8 [2]	18.6 [3]	9.6 [1.6]	6.1 [2.9]	14.1 [3.8]	7.4 [2.1]
<i>scatt-NMF<sub>1</sub></i>	10.3 [2.0]	19.7 [3.3]	11.0 [1.7]	6.2 [2.8]	13.5 [3.5]	7.8 [2.2]
<i>scatt-NMF<sub>2</sub></i>	<b>10.6</b> [1.8]	<b>20.5</b> [3.0]	<b>11.3</b> [1.7]	<b>6.9</b> [2.7]	<b>16.0</b> [3.5]	<b>7.9</b> [2.2]

**Table 1:** Separation with speakers-specific and multi-speaker settings. Average SDR, SIR and SAR (in *dB*) for NMF and proposed and *scatt-NMF<sub>2</sub>*. Standard deviation of each result shown between brackets.

**Data sets.** We used a subset of the GRID dataset [24] for evaluating the speaker-specific setting. For each speaker, 500 randomly-chosen clips were used for training (around 25 minutes) and 200 clips were used for testing. For the multi-speaker case we used a subset of the TIMIT dataset. We adopted the standard test-train division, using all the training recordings for building the models and a subset of 12 different speakers (6 males and 6 females) for testing. For each speaker we randomly chose two clips and compared all female-male combinations (144 mixtures).

**Evaluation measures.** We used the *source-to-distortion ratio* (SDR), *source-to-interference ratio* (SIR), and *source-to-artifact ratio* (SAR) from the BSS-EVAL metrics [25].

**Training setting.** We evaluated the proposed scattering NMF model with one and two layers, referred as *scatt-NMF<sub>1</sub>* and *scatt-NMF<sub>2</sub>* respectively. As a baseline we used standard NMF with frame lengths of 1024 samples and 50% overlap. The dictionaries in standard NMF were chosen with 200 and 400 atoms for the speaker-specific and multi-speaker settings respectively. These values were obtained using cross-validation on a few clips separated from the training as a validation set. In all cases, we applied *scatt-NMF* using a scattering transforms with resolution  $Q_1 = 32$  and  $Q_2 = 1$ . The resulting representation had 175 coefficients for the first level and around 2000 for the second layer. For the single speaker case we trained dictionaries with 200 atoms for *scatt-NMF<sub>1</sub>* and 800 atoms for *scatt-NMF<sub>2</sub>*. While for the multi-speaker case we used 400 atoms for *scatt-NMF<sub>1</sub>* and 1000 atoms for *scatt-NMF<sub>2</sub>*. In all cases, the features were frame-wise normalized and we used  $\lambda = 0.1$ .

**Results.** Table 1 shows the results obtained for the speaker-specific and multi-speaker settings.<sup>1</sup> In all cases we observe that the one layer scattering transform outperforms the STFT in terms of SDR. Furthermore, there is a tangible gain in including a deeper representation; *scatt-NMF<sub>2</sub>* performs always better than *scatt-NMF<sub>1</sub>*. As expected, the results obtained with the speaker-specific setting are better than those of the more challenging problem of multi-speaker setting.

We also compared the proposed approach with the speaker-specific setting discussed in [26]. In this work the authors investigate several alternatives of using Recurrent Neural Networks (RNN) for speech separation. Several optimization

	SDR	SIR	SAR
NMF-KL	5.4	7.3	7.8
RNN [26]	6.0	8.1	<b>8.1</b>
RNN joint disc. training [26]	<b>7.4</b>	<b>11.8</b>	7.5
<i>scatt-NMF<sub>2</sub></i>	6.7	11.1	6.9

**Table 2:** Comparison RNN based separation, with and without joint discriminative training with soft masks [26].

settings are evaluated on two given speakers of the TIMIT dataset, some of which aim at learning short-term temporal dynamics. This is a very challenging setting due to the very small available training data (less than 10 seconds per speaker). The evaluations of *scatt-NMF<sub>2</sub>* were performed using the setting provided in [26] (with the corresponding training, development and testing data) while their results are taken from the paper. *scatt-NMF<sub>2</sub>* outperforms the benchmark KL-NMF in SDR and SIR, and is competitive with the best performing networks reported in [26], with and without joint discriminative training, see Table 2. We expect further gains by applying discriminative dictionary learning [27].

In summary, these results confirm that inverse problems such as speech source separation can benefit from the properties of stable and highly discriminative non-linear representations, such as scattering operators. Sparse inference is able to extract more relevant information thanks to the stability to time-frequency deformations, while the phase recovery can still be efficiently performed with gradient descent.

## 6. DISCUSSION

NMF-based audio source separation techniques can be thought as applying a synthesis operator on a feature space given by a pooled analysis operator. Leveraging recent developments in signal processing, we propose to substitute the first stage with a deep scattering transform. The obtained features are designed to capture the joint time-frequency variability of speech signals and efficiently represent a longer temporal context. Experimental evaluation shows that using deeper representations leads to a tangible improvement in performance in challenging source separation settings. A natural extension of this work is to investigate the use of learned representations instead, or on top of, the designed ones. Future work includes testing more thoroughly the potential of the proposed model in combination with convolutional neural networks, which have been very successful in other signal and image processing problems.

<sup>1</sup>Audio samples are available at [www.cims.nyu.edu/~bruna/scatt\\_source\\_separation](http://www.cims.nyu.edu/~bruna/scatt_source_separation).

## 7. REFERENCES

- [1] P. C. Loizou, *Speech Enhancement: Theory and Practice*, vol. 30, CRC, 2007.
- [2] E. Hänsler and G. Schmidt, *Speech and Audio Processing in Adverse Environments*, Springer, 2008.
- [3] D.D. Lee and H.S. Seung, “Learning parts of objects by non-negative matrix factorization,” *Nature*, vol. 401, no. 6755, pp. 788–791, 1999.
- [4] P. Smaragdis, C. Févotte, G. Mysore, N. Mohammadiha, and M. Hoffman, “Static and dynamic source separation using nonnegative factorizations: A unified view,” *Signal Processing Magazine, IEEE*, vol. 31, no. 3, pp. 66–75, 2014.
- [5] M. N. Schmidt and R. K. Olsson, “Single-channel speech separation using sparse non-negative matrix factorization,” in *INTERSPEECH*, Sep 2006.
- [6] M. V. S. Shashanka, B. Raj, and P. Smaragdis, “Sparse Overcomplete Decomposition for Single Channel Speaker Separation,” in *ICASSP*, 2007.
- [7] Z. Duan, G. J. Mysore, and P. Smaragdis, “Online plca for real-time semi-supervised source separation,” in *LVA/ICA*, 2012, pp. 34–41.
- [8] N. Mohammadiha, P. Smaragdis, and A. Leijon, “Supervised and unsupervised speech enhancement using nonnegative matrix factorization,” *Audio, Speech, and Language Processing, IEEE Transactions on*, vol. 21, no. 10, pp. 2140–2151, 2013.
- [9] C. Févotte, “Majorization-minimization algorithm for smooth itakura-saito nonnegative matrix factorization,” in *ICASSP. IEEE*, 2011, pp. 1980–1983.
- [10] K. W. Wilson, B. Raj, P. Smaragdis, and A. Divakaran, “Speech denoising using nonnegative matrix factorization with priors,” in *ICASSP*, 2008, pp. 4029–4032.
- [11] G. J. Mysore and P. Smaragdis, “A non-negative approach to semi-supervised separation of speech from noise with the use of temporal dynamics,” in *ICASSP*, 2011, pp. 17–20.
- [12] J. Han, G. J. Mysore, and B. Pardo, “Audio imputation using the non-negative hidden markov model,” in *LVA/ICA*, 2012, pp. 347–355.
- [13] C. Févotte, J. Le Roux, and J. R. Hershey, “Non-negative dynamical system with application to speech and audio,” in *ICASSP*, 2013.
- [14] R. W. Gerchberg and W. Owen Saxton, “A practical algorithm for the determination of the phase from image and diffraction plane pictures,” *Optik*, vol. 35, pp. 237–246, 1972.
- [15] J. Andén and S. Mallat, “Deep scattering spectrum,” *arXiv preprint arXiv:1304.6763*, 2013.
- [16] J. Bruna and S. Mallat, “Invariant scattering convolution networks,” *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 35, no. 8, pp. 1872–1886, 2013.
- [17] J. Bruna, *Scattering Representations for Recognition*, Ph.D. thesis, Palaiseau, Ecole polytechnique, 2013.
- [18] G. Wolf, S. Mallat, and S. Shamma, “Audio source separation with time-frequency velocities,” *International Workshop on Machine Learning for Signal Processing*, 2014.
- [19] J. Bruna and S. Mallat, “Audio texture synthesis with scattering moments,” *arXiv preprint arXiv:1311.0407*, 2013.
- [20] R. Jenatton, R. Gribonval, and F. Bach, “Local stability and robustness of sparse dictionary learning in the presence of noise,” *arXiv preprint arXiv:1210.0685*, 2012.
- [21] C. Févotte and J. Idier, “Algorithms for nonnegative matrix factorization with the  $\beta$ -divergence,” *Neural Computation*, vol. 23, no. 9, pp. 2421–2456, 2011.
- [22] M. N. Schmidt, J. Larsen, and F.-T. Hsiao, “Wind noise reduction using non-negative sparse coding,” in *MLSP*, Aug 2007, pp. 431–436.
- [23] J. Mairal, F. Bach, J. Ponce, and G. Sapiro, “Online dictionary learning for sparse coding,” in *ICML*, 2009, pp. 689–696.
- [24] M. Cooke, J. Barker, S. Cunningham, and X. Shao, “An audio-visual corpus for speech perception and automatic speech recognition,” *J. of the Acoustical Society of America*, vol. 120, pp. 2421, 2006.
- [25] E. Vincent, R. Gribonval, and C. Févotte, “Performance measurement in blind audio source separation,” *IEEE Trans. on Audio, Speech, and Lang. Proc.*, vol. 14, no. 4, pp. 1462–1469, 2006.
- [26] P.-S. Huang, M. Kim, M. Hasegawa-Johnson, and P. Smaragdis, “Deep learning for monaural speech separation,” in *ICASSP*, 2014, pp. 1562–1566.
- [27] P. Sprechmann, A. M. Bronstein, and G. Sapiro, “Supervised non-euclidean sparse NMF via bilevel optimization with applications to speech enhancement,” in *HSCMA. IEEE*, 2014, pp. 11–15.