

*This is the peer reviewed version of the following article: [Juan De La Torre Cruz, Francisco Jesús Cañadas Quesada, Damián Martínez-Muñoz, Nicolás Ruiz Reyes, Sebastián García Galán, Julio José Carabias Orti. An incremental algorithm based on multichannel non-negative matrix partial co-factorization for ambient denoising in auscultation. Applied Acoustics Volume 182, November 2021, 108229], which has been published in final form at <https://doi.org/10.1016/j.apacoust.2021.108229>] This article may be used for non-commercial purposes in accordance with Applied Acoustics and conditions for use of self-archived versions. This article may not be enhanced, enriched or otherwise transformed into a derivative work, without express permission from Applied Acoustics or by statutory rights under applicable legislation. Copyright notices must not be removed, obscured or modified. The article must be linked to Applied Acoustics's version of record on Applied Acoustics online library and any embedding, framing or otherwise making available the article or pages thereof by third parties from platforms, services and websites other than Applied Acoustics online library must be prohibited.*

# An incremental algorithm based on multichannel non-negative matrix partial co-factorization for ambient denoising in auscultation

Juan De La Torre Cruz \*, Francisco Jesús Cañadas Quesada, Damián Martínez-Muñoz, Nicolás Ruiz Reyes, Sebastián García Galán, Julio José Carabias Orti

*Department of Telecommunication Engineering. University of Jaen, Campus Científico-Tecnológico de Linares, Avda. de la Universidad, s/n, 23700 Linares, Jaen, Spain*

---

## Abstract

One of the major current limitations in the diagnosis derived from auscultation remains the ambient noise surrounding the subject, which prevents successful auscultation. Therefore, it is essential to develop robust signal processing algorithms that can extract relevant clinical information from auscultated recordings analyzing in depth the acoustic environment in order to help the decision-making process made by physicians. The aim of this study is to implement a method to remove ambient noise in biomedical sounds captured in auscultation. We propose an incremental approach based on multichannel non-negative matrix partial co-factorization (NMPCF) for ambient denoising focusing on high noisy environment with a Signal-to-Noise Ratio (SNR)  $\leq -5$  dB. The first contribution applies NMPCF assuming that ambient noise can be modelled as repetitive sound events simultaneously found in two single-channel inputs captured by means of different recording devices. The second contribution proposes an incremental algorithm, based on the previous multichannel NMPCF, that refines the estimated biomedical spectrogram throughout a set of incremental stages

---

\*Corresponding author. Tlf.: (+34) 953648592

*Email addresses:* [jtorre@ujaen.es](mailto:jtorre@ujaen.es) (Juan De La Torre Cruz \*), [fcandas@ujaen.es](mailto:fcandas@ujaen.es) (Francisco Jesús Cañadas Quesada), [damian@ujaen.es](mailto:damian@ujaen.es) (Damián Martínez-Muñoz), [nicolas@ujaen.es](mailto:nicolas@ujaen.es) (Nicolás Ruiz Reyes), [sgalan@ujaen.es](mailto:sgalan@ujaen.es) (Sebastián García Galán), [carabias@ujaen.es](mailto:carabias@ujaen.es) (Julio José Carabias Orti)

by eliminating most of the ambient noise that was not removed in the previous stage at the expense of preserving most of the biomedical spectral content. The ambient denoising performance of the proposed method, compared to some of the most relevant state-of-the-art methods, has been evaluated using a set of recordings composed of biomedical sounds mixed with ambient noise that typically surrounds a medical consultation room to simulate high noisy environments with a SNR from -20 dB to -5 dB. In order to analyse the drop in denoising performance of the evaluated methods when the effect of the propagation of the patient's body material and the acoustics of the room is considered, results have been obtained with and without taking these effects into account. Experimental results report that: (i) the performance drop suffered by the proposed method is lower compared to MSS and NLMS when considering the effect of the propagation of the patient's body material and the acoustics of the room active; (ii) unlike what happens with MSS and NLMS, the proposed method shows a stable trend of the average SDR and SIR results regardless of the type of ambient noise and the SNR level evaluated; and (iii) a remarkable advantage of the proposed method is the high robustness of the acoustic quality of the estimated biomedical sounds when the two single-channel inputs suffer from a delay between them.

*Keywords:* Auscultation, Biomedical, Ambient noise, Non-negative matrix partial co-factorization, Multichannel, Incremental

---

## 1. Introduction

Auscultation is defined as the technique of listening the internal sounds produced by the human organs by means of a stethoscope. This technique is simple, non-invasive, safe and inexpensive that provides valuable clinical information in the diagnosis of the status of the heart, lung and airways [1, 2]. Although today there are more advanced medical tools to analyze the status of the heart and lung, such as chest radiography, electrocardiography (ECG), spirometry or laboratory analyses, auscultation is still one of the most widely used techniques

9 to detect any cardiac or pulmonary disease. However, the diagnosis derived  
10 from auscultation shows two main limitations: i) high subjectivity due to the  
11 physician's expertise to recognize sounds that reveal any physiological disorder  
12 [3]; ii) high dependence on the ambient noise surrounding the subject to provide  
13 a reliable diagnosis [4].

14 Because the process of auscultation in a soundproof room is not possible in  
15 most cases, especially in low-income and middle-income countries supported by  
16 a resource-poor health system, ambient denoising performed in the examination  
17 room of the health center is still a challenging task in biomedical signal process-  
18 ing in order to maximize the reliability of a diagnosis. The main effects caused  
19 by ambient noise are the masking, distortion and weakness of the sound of in-  
20 terest that may provide relevant clues in the diagnosis as shown in Figure 1. As  
21 a result, the probability of making a medical error increases when auscultation  
22 is performed in a noisy environment since the physician is not able to correctly  
23 interpret the diagnostic information contained in the sound signal from auscul-  
24 tation. In this work, the term biomedical means that the sound sources that  
25 have generated the sounds of interest have been the human internal organs, and  
26 specifically, the heart and the lung.

27 In recent years, several signal processing tasks have been applied in the field  
28 of biomedical information retrieval such as, sound source separation [5, 6, 7] as  
29 well as sound event detection [8, 9, 10, 11] and classification [12, 13, 14, 15, 16].  
30 However, most of their experimental results have been obtained in environments  
31 in which the biomedical sounds are not acoustically contaminated by ambient  
32 noises. Therefore, the task of ambient denoising is still an open research topic  
33 in biomedical engineering being most of the approaches based on adaptive fil-  
34 tering [17, 18, 19, 20, 21] and spectral subtraction [22, 23, 24, 25]. Chang  
35 and Lai [23] proposed a two-channel spectral subtraction method, based on  
36 autoregressive (AR), mel-frequency cepstral coefficients (MFCC) and dynamic  
37 time warping (DTW), applied to the lung sound signals under noisy conditions  
38 before the extraction of lung sound features. Emmanouilidou et al. [24] devel-  
39 oped a multiband spectral scheme, based on two-microphone setup, to suppress

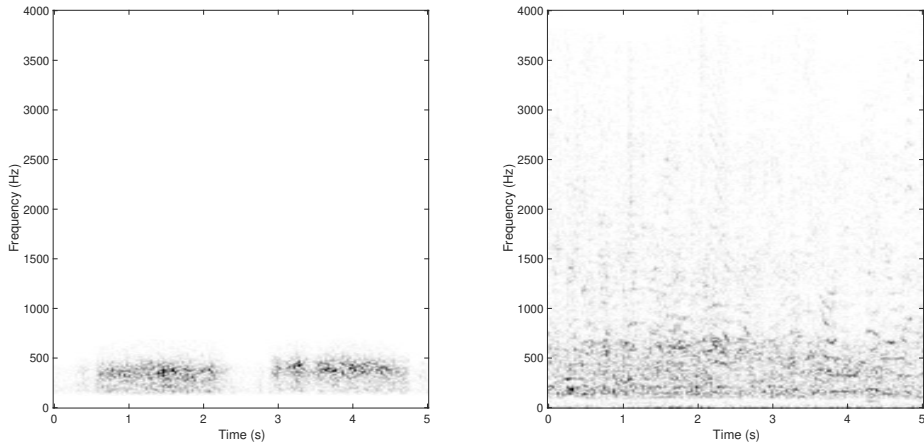


Figure 1: Spectrogram from a clean lung sound recording: (Left) under no ambient noise; (Right) mixed with ambient noise (babble) in a Signal-to-Noise Ratio (SNR) equals -10 dB. Comparing both figures, it can be observed that the spectral content from lung is completely masked by ambient noise. Higher energies are indicated by darker colour.

40 the background noise while successfully preserving the lung sound content to  
 41 maximize the informative diagnostic value obtained from auscultation. The al-  
 42 gorithm analyzes each frequency band in a nonuniform manner and uses prior  
 43 knowledge of the target sounds to apply a penalty in the spectral domain. It  
 44 follows from the above that it is crucial to develop robust signal processing algo-  
 45 rithms that can extract relevant clinical information from auscultated recordings  
 46 taking into consideration the acoustic environment that surrounds the subject  
 47 in order to improve the decision making process made by physicians.

48 It is well known that the conventional Non-negative Matrix Partial Co-  
 49 Factorization (NMPCF) enforces a joint matrix decomposition using multiple  
 50 matrices to recover a set of shared spectral patterns (bases) that model the  
 51 spectral behavior of some of the sound sources contained in the single-channel  
 52 input. Over the last decade, NMPCF has been successfully applied in several  
 53 single-channel sound signal processing fields: i) Music Information Retrieval  
 54 (MIR) such as, singing-voice separation [26], rhythmic extraction [27, 28, 29]  
 55 and speaker diarization [30]; and ii) Enhancement of biomedical sounds such  
 56 as, normal respiratory and wheezes sound separation [31]. In this work, we pro-

57 pose an incremental algorithm, called 2C-NMPCF, that improves the quality  
58 of biomedical sounds captured in auscultated recordings by applying the con-  
59 ventional NMPCF from a multi-channel scenario rather than a single-channel.  
60 In this paper, the term multichannel refers to the use of two single-channel  
61 audio inputs simultaneously captured by means of different recording devices.  
62 As occurs in [24], these two single-channel inputs are defined as: (i) the inter-  
63 nal recording that comes from the audio captured using a stethoscope in which  
64 both biomedical sounds from inside the human body and ambient noises can be  
65 listened; (ii) the external recording that comes from the audio captured using  
66 an external microphone in which only the ambient noise that surrounds the  
67 subject is captured. Specifically, our first contribution applies NMPCF from  
68 a multichannel point of view assuming that ambient noises can be modelled  
69 as repetitive sounds that can be simultaneously found in both single-channel  
70 inputs. In other words, we implicitly assume that the spectral patterns that  
71 characterize the ambient noises are repeated sound events contained in both  
72 the spectrograms from the internal and external recordings. Our second contri-  
73 bution proposes an incremental algorithm, based on the previous multichannel  
74 NMPCF, that refines the estimated biomedical spectrogram through a set of  
75 incremental stages by eliminating a high amount of ambient noise that was not  
76 extracted in the previous stage, especially in the case of high noise environ-  
77 ments. In this work, a high noisy environment provides a Signal-to-Noise Ratio  
78 (SNR) lower than 0 dB.

79 The paper is structured as follows: Section 2 details the datasets, the state-  
80 of-the-art methods for comparison and the proposed method. The metrics, setup  
81 and results are shown in Section 3. Finally, Section 4 presents the conclusions  
82 and future work.

## 83 2. Materials and Methods

### 84 2.1. Data collection

85 Due to the lack of publicly available databases consisting of biomedical  
86 sounds mixed with ambient noises to the best of our knowledge, we have cre-  
87 ated the database  $D_C$ . The database  $D_C$  is composed by the ambient noise  
88 database  $D_N$  and the biomedical database  $D_B$  in order to simulate auscultation  
89 recordings captured from a stethoscope.

90 The database  $D_N$  has been created taking into account a wide range of  
91 ambient noises collected from databases widely used in the field of sound source  
92 separation [32] and sound event detection [33, 34]. Most of these ambient noises  
93 have been classified as some of the most disturbing noises that can appear in the  
94 auscultation performed in the hospital room according to information provided  
95 by medical personnel from the Hospital of Jaen (Spain). For this reason, the  
96 database  $D_N$  is composed of five types of ambient noise in order to assess the  
97 denoising performance of the proposed method considering common indoor and  
98 outdoor ambient noises that typically surround a medical consultation room:  
99 ambulance siren [35, 36], baby crying [37], babble (people speaking) [38, 39],  
100 car (inside the vehicle) [40] and street (car passing by, car engine running, car  
101 idling, bus, truck, children yelling, people talking, workers on the street) [41, 42].  
102 The database  $D_N$  consists of a total of 150 single-channel recordings of ambient  
103 noises, of which each type of noise consists of 30 recordings. Each recording  
104 has a duration of 5 seconds that has been obtained applying a pseudo-random  
105 process, based on the standard uniform distribution, to select a starting time  
106 followed by a 5 seconds interval.

107 The database  $D_B$  consists of a total of 150 single-channel biomedical record-  
108 ings from public and private biomedical databases, specifically, 75 heart record-  
109 ings [43, 44] (typically in the range 10Hz-320 Hz [45, 46]) and 75 lung recordings  
110 [47] (typically in the range 50Hz-2500 Hz [48, 49]). Highlight that ambient noises  
111 are not listened on each recording. Each recording has a duration of 5 seconds  
112 which has been obtained applying a pseudo-random process similarly as used in

113 the database  $D_N$ .

114 Each mixture recording belonging to the database  $D_C$  has been generated  
115 mixing each recording from the database  $D_B$  with a recording of each type of  
116 noise randomly chosen from the database  $D_N$ . Indicate that the recordings of  
117 noise used for the mixtures with the heart recordings are the same as those  
118 used for the mixing with the lung recordings. For each mixture recording from  
119  $D_C$ , the ambient noise used in the internal recording is the same noise used  
120 in the external recording. The database  $D_C$  is not affected by the effect of  
121 the patient's body material or by the acoustics room in order to perform a  
122 fair optimization parameters of the proposed method avoiding body-specific or  
123 room-specific optimization. As a result, two databases are created from  $D_C$ ,  
124 the optimization database  $D_O$  and the testing database  $D_T$ ,

- 125 • The optimization database  $D_O$  is generated randomly selecting two-thirds  
126 of all mixtures recordings from the database  $D_C$ .
- 127 • The testing database  $D_T$  is generated using the remainder one-third mix-  
128 tures recordings that are not used in the database  $D_O$ .

129 The set of recordings used in the optimization database  $D_O$  is not the same  
130 as that used in the testing database  $D_T$  in order to validate the denoising results.  
131 Moreover, several SNR have been applied in the mixing process to create the  
132 database  $D_C$  in order to evaluate high noisy environments. In this way, the  
133 databases  $D_{T_{-20}}$  (SNR=-20 dB),  $D_{T_{-15}}$  (SNR=-15 dB),  $D_{T_{-10}}$  (SNR=-10 dB)  
134 and  $D_{T_{-5}}$  (SNR=-5 dB) refer to the same database  $D_T$  but using different  
135 SNR between biomedical and ambient noise recordings. For example, a value  
136 SNR=-20 dB indicates that the power of the ambient noise is 100 times greater  
137 compared to the power of the biomedical sounds used in the audio mixture.  
138 Table 1 describes the characteristics of the data, according to the databases  
139 used in the experimental evaluation.

$ID_1$	$ID_2$	$ID_3$	$ID_4$	$ID_5$	$ID_6$	$ID_7$	$ID_8$	$ID_9$	$ID_{10}$	$ID_{11}$
$D_B$	75	75	-	-	-	-	-	150	-	750
$D_N$	-	-	30	30	30	30	30	-	150	750
$D_O$	50	50	20	20	20	20	20	100	100	2500
$D_T$	25	25	10	10	10	10	10	50	50	1250

Table 1: Characteristics of the data.  $ID_1$ : database identifier;  $ID_2$ : number of clean heart recordings;  $ID_3$ : number of clean lung recordings;  $ID_4$ : number of ambulance siren noise recordings;  $ID_5$ : number of babycriing noise recordings;  $ID_6$ : number of babble noise recordings;  $ID_7$ : number of car noise recordings;  $ID_8$ : number of street noise recordings;  $ID_9$ : total number of biomedical (clean heart and lung) recordings;  $ID_{10}$ : total number of ambient noise recordings;  $ID_{11}$ : temporal duration for all recordings in seconds.

140 *2.2. Baseline method for comparison*

141 A reference adaptive filtering method, Normalized Least-Mean-Square (NLMS)  
142 [50, 51], has been used to assess the performance of the proposed method in the  
143 task of ambient denoising. The optimal length of the filter has been optimized  
144 by varying the number of coefficients between 10 and 100 to maximize the de-  
145 noising performance of the NLMS method. The optimization results showed  
146 that NLMS obtained the best denoising performance when the size of the adap-  
147 tive filter is equal to 10. Moreover, one of the most relevant state-of-the-art  
148 methods, based on Multiband Spectral Subtraction (MSS) [24], has also been  
149 implemented. Note that the best configuration (algorithm B) of MSS has been  
150 chosen to perform a fair comparison. The reader can refer to [24] for more  
151 details.

152 *2.3. Proposed method for ambient denoising*

153 The main problem that physicians point out when performing the auscultation  
154 process in high noisy environments is that the biomedical sounds are  
155 severely overlapped with ambient noises so, part of the valuable clinical infor-  
156 mation contained in the sounds of interest is masked. The aim of the pro-  
157 posed method is to improve the quality of the biomedical sounds captured by  
158 a stethoscope in high noisy environments applying Non-negative Matrix Partial



159 Co-Factorization (NMPCF) in a multichannel (two distinct single-channel sig-  
 160 nals) scenario (2C-NMPCF). The flowchart of the proposed method 2C-NMPCF  
 161 is shown in Figure 2.

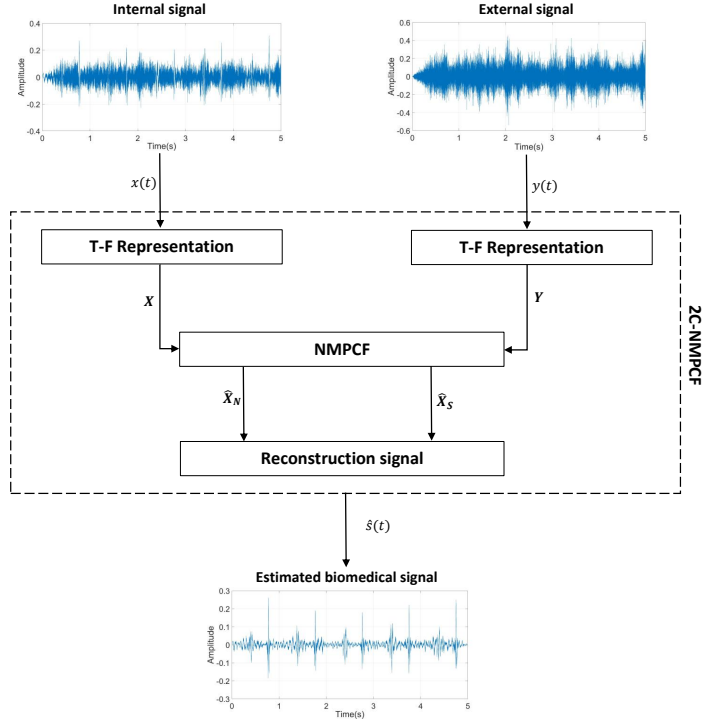


Figure 2: Flowchart of the proposed method 2C-NMPCF.

162 *2.3.1. Time-Frequency representation*

163 The internal signal (first single-channel)  $x(t)$  represents the sounds captured  
 164 by a digital stethoscope that is composed of two types of additive sound sources:  
 165 (i) the biomedical sounds  $s(t)$  from the subject; (ii) the ambient noises  $n(t)$   
 166 surrounding the subject that are still heard inside the human body. We assume  
 167 that  $s(t)$  and  $n(t)$  can be considered independent sound sources, that is,  $x(t) =$   
 168  $s(t) + n(t)$ . Moreover, an external microphone, located outside of the subject,  
 169 captures all ambient noises that are represented by the external signal (second  
 170 single-channel)  $y(t)$ .

171 The complex and magnitude spectrogram  $\mathbf{X}_c \in \mathbb{C}_+^{F \times T}$ ,  $\mathbf{X} \in \mathbb{R}_+^{F \times T}$  asso-  
172 ciated to the internal signal  $x(t)$  and the magnitude spectrogram  $\mathbf{Y} \in \mathbb{R}_+^{F \times T}$   
173 associated to the external signal  $y(t)$  are calculated using the Short-Time Fourier  
174 Transform (STFT) applying a Hamming window of size  $N$  with 50% overlap.  
175 Indicate that the complex values associated with  $\mathbf{X}_c$ , in which the phase in-  
176 formation is included, are used later in the resynthesis process. The size and  
177 scale of the magnitude spectrograms depend on each input single-channel sig-  
178 nal. Therefore, a normalization process is applied in order to ensure that the  
179 proposed method is independent of the size and scale of the input spectrograms.  
180 Thus, the normalized spectrograms  $\bar{\mathbf{X}} \in \mathbb{R}_+^{F \times T}$  and  $\bar{\mathbf{Y}} \in \mathbb{R}_+^{F \times T}$  are computed  
181 as follows,

$$\bar{\mathbf{Z}} = \frac{\mathbf{Z}}{\left( \frac{\sum_{f,t} Z_{f,t}}{FT} \right)} \quad (1)$$

182 where  $\mathbf{Z} = \{\mathbf{X}, \mathbf{Y}\}$  according to the input spectrogram. The variables  $F$  and  
183  $T$  represent the number of frequency bins and the number of time frames. To  
184 avoid the complex nomenclature throughout the manuscript, the variables  $\bar{\mathbf{X}}$   
185 and  $\bar{\mathbf{Y}}$  are hereinafter referred as  $\mathbf{X}$  and  $\mathbf{Y}$ , respectively.

### 186 2.3.2. Multichannel Non-negative Matrix Partial Co-Factorization (2C-NMPCF)

187 The idea of the proposed method 2C-NMPCF is to enforce a joint matrix de-  
188 composition using multiple matrices  $\mathbf{X}$ ,  $\mathbf{Y}$  obtained from distinct single-channel  
189 spectrograms instead of the several excerpts of the same single-channel spectro-  
190 gram as occurs in the conventional NMPCF. The main contribution of the pro-  
191 posed method 2C-NMPCF is to exploit the spectral patterns that are shared in  
192 two distinct spectrograms since we assume that ambient noises can be modelled  
193 as repetitive sound events that can be simultaneously found in the spectrograms  
194 associated both the internal and external signal. This modeling allows to remove  
195 most of the ambient noises that are active in the internal signal improving the  
196 quality of the biomedical sounds from the auscultation process. The proposed  
197 method 2C-NMPCF is composed of two stages:

- 198 • Stage 1. This stage is applied to the internal signal  $x(t)$ . The input

199 spectrogram  $\mathbf{X}$  is decomposed into two separated or estimated spectro-  
 200 grams, the magnitude spectrogram only composed of biomedical sounds  
 201  $\hat{\mathbf{X}}_S \in \mathbb{R}_+^{F \times T}$  and the magnitude spectrogram only composed of ambient  
 202 noises  $\hat{\mathbf{X}}_N \in \mathbb{R}_+^{F \times T}$ . The factorization of each spectrogram depends on the  
 203 estimated basis matrix  $\mathbf{U} \in \mathbb{R}_+^{F \times K}$  (dictionary of spectral patterns) and  
 204 the estimated activation matrix  $\mathbf{V} \in \mathbb{R}_+^{K \times T}$  (temporal gains) as follows,

$$\mathbf{X} \approx \hat{\mathbf{X}} = \hat{\mathbf{X}}_N + \hat{\mathbf{X}}_S = \mathbf{U}\mathbf{V} = \begin{bmatrix} \mathbf{U}_N & \mathbf{U}_S \end{bmatrix} \begin{bmatrix} \mathbf{V}_N \\ \mathbf{V}_S \end{bmatrix} = \mathbf{U}_N\mathbf{V}_N + \mathbf{U}_S\mathbf{V}_S \quad (2)$$

205 where  $\hat{\mathbf{X}} \in \mathbb{R}_+^{F \times T}$  is the estimated or reconstructed magnitude spectro-  
 206 gram of the first channel signal.  $\mathbf{U}_N \in \mathbb{R}_+^{F \times K_N}$  and  $\mathbf{V}_N \in \mathbb{R}_+^{K_N \times T}$  are  
 207 the estimated basis and activations matrix of the ambient noises. The  
 208 variables  $\mathbf{U}_S \in \mathbb{R}_+^{F \times K_S}$  and  $\mathbf{V}_S \in \mathbb{R}_+^{K_S \times T}$  are the estimated basis and ac-  
 209 tivation matrix of the biomedical sounds. The parameter  $K = K_N + K_S$   
 210 indicates the number of bases, being  $K_N$  the number of bases related to  
 211 the ambient noises and  $K_S$  the number of bases related to the biomedical  
 212 sounds. In this stage, the described decomposition model (see Equation  
 213 (3)) does not obtain a parts-based objects reconstruction with physical  
 214 meaning as occurs in real-world. Therefore, this stage cannot distinguish  
 215 between spectral patterns belonging to biomedical sounds and ambient  
 216 noises.

- 217 • Stage 2. This stage is applied to the external signal  $y(t)$ . We assume  
 218 that the external signal is only composed of ambient noises, therefore the  
 219 goal of this model is to reconstruct the external magnitude spectrogram  
 220  $\mathbf{Y}$  by using the basis matrix  $\mathbf{U}_N$  composed of the spectral patterns that  
 221 characterize the ambient noises,

$$\mathbf{Y} \approx \hat{\mathbf{Y}} = \mathbf{U}_N\mathbf{H}_N \quad (3)$$

222 where  $\hat{\mathbf{Y}} \in \mathbb{R}_+^{F \times T}$  is the estimated or reconstructed magnitude spectro-  
 223 gram of the external signal. The variable  $\mathbf{H}_N \in \mathbb{R}_+^{K_N \times T}$  is the estimated

224 activations matrix of the ambient noises for the external signal. Note that  
 225  $\mathbf{U}_N$  can be treated as the same matrix previously used in Equation (2).

226 Specifically, 2C-NMPCF extends the conventional NMPCF to a multichan-  
 227 nel scenario sharing the frequency basis matrix  $\mathbf{U}_N$  to factorize simultaneously  
 228 two distinct single-channel magnitude spectrograms  $\mathbf{X}$ ,  $\mathbf{Y}$  as shown in Figure 3.  
 229 The proposed method 2C-NMPCF allows to factorize jointly  $\mathbf{X}$  and  $\mathbf{Y}$  so the  
 230 spectral patterns of the ambient noises, active in both spectrograms, are shared  
 231 in the same dictionary  $\mathbf{U}_N$  since we assume that ambient noises can be con-  
 232 sidered repetitive sounds that can be simultaneously active in both magnitude  
 233 spectrograms  $\mathbf{X}$  and  $\mathbf{Y}$ . Contrarily, the dictionary  $\mathbf{U}_S$  represents the spec-  
 234 tral patterns of the biomedical sounds that only can be found in the internal  
 235 magnitude spectrogram  $\mathbf{X}$ .

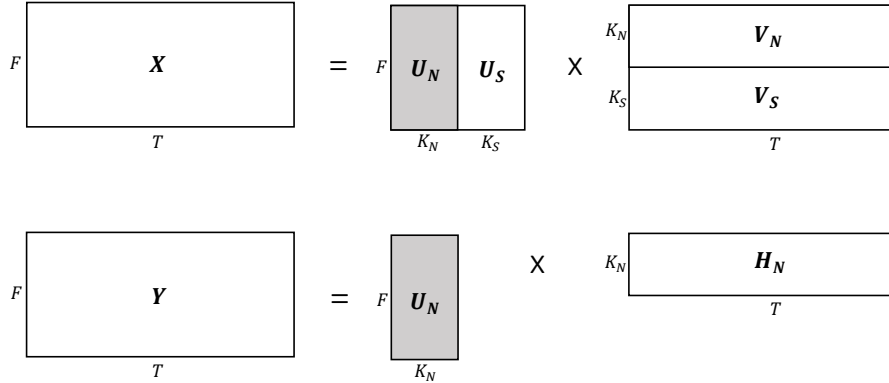


Figure 3: Matrix decomposition based on multichannel NMPCF (2C-NMPCF).

### 236 2.3.3. Objective Function and Update rules

237 The objective function of the proposed method 2C-NMPCF that must be  
 238 performed to minimize the residuals of the two previous models, see Equations  
 239 (2)-(3), is detailed as follows,

$$\Gamma = D_{KL}(\mathbf{X}|\hat{\mathbf{X}}) + \lambda D_{KL}(\mathbf{Y}|\hat{\mathbf{Y}}) \quad (4)$$

240 where the parameter  $\lambda$  controls the relative importance between the internal  
 241 and the external magnitude spectrogram. So, the contribution of the magnitude  
 242 spectrogram  $\mathbf{Y}$  is greater when the parameter  $\lambda$  increases. In this paper, the  
 243 Kullback–Leibler divergence, see Equation (5), has been used to calculate the  
 244 signal reconstruction error for the internal spectrogram  $D_{KL}(\mathbf{X}|\hat{\mathbf{X}})$  and the  
 245 external spectrogram  $D_{KL}(\mathbf{Y}|\hat{\mathbf{Y}})$ . The reason is because this cost function  
 246  $D_{KL}$  is non-increasing, ensuring the non-negativity of the estimated basis and  
 247 activations matrices and moreover, several works have demonstrated promising  
 248 results in the field of biomedical signal processing [7, 11, 52].

$$D_{KL}(\mathbf{Z}|\hat{\mathbf{Z}}) = \mathbf{Z} \log \frac{\mathbf{Z}}{\hat{\mathbf{Z}}} - \mathbf{Z} + \hat{\mathbf{Z}}, \quad \mathbf{Z} = \{\mathbf{X}, \mathbf{Y}\} \quad (5)$$

249 From Equation (4), the estimated basis matrices  $\mathbf{U}_N, \mathbf{U}_S$  and activation  
 250 matrices  $\mathbf{V}_N, \mathbf{V}_S, \mathbf{H}_N$  can be obtained by applying a gradient descent algorithm  
 251 based on multiplicative update rules. The multiplicative update rules to learn  
 252 those matrices can be obtained by taking negative and positive terms of the  
 253 partial derivative of the cost function  $\Gamma$  with respect to  $\mathbf{U}_N, \mathbf{U}_S, \mathbf{V}_N, \mathbf{V}_S$  and  
 254  $\mathbf{H}_N$ , respectively,

$$\mathbf{U}_N \leftarrow \mathbf{U}_N \odot \frac{(\mathbf{X} \oslash \hat{\mathbf{X}})(\mathbf{V}_N)^T + \lambda(\mathbf{Y} \oslash \hat{\mathbf{Y}})(\mathbf{H}_N)^T}{(\mathbf{V}_N)^T + \lambda(\mathbf{H}_N)^T} \quad (6)$$

$$\mathbf{U}_S \leftarrow \mathbf{U}_S \odot \frac{(\mathbf{X} \oslash \hat{\mathbf{X}})(\mathbf{V}_S)^T}{(\mathbf{V}_S)^T} \quad (7)$$

$$\mathbf{V}_N \leftarrow \mathbf{V}_N \odot \frac{(\mathbf{U}_N)^T(\mathbf{X} \oslash \hat{\mathbf{X}})}{(\mathbf{U}_N)^T} \quad (8)$$

$$\mathbf{V}_S \leftarrow \mathbf{V}_S \odot \frac{(\mathbf{U}_S)^T(\mathbf{X} \oslash \hat{\mathbf{X}})}{(\mathbf{U}_S)^T} \quad (9)$$

$$\mathbf{H}_N \leftarrow \mathbf{H}_N \odot \frac{(\mathbf{U}_N)^T(\mathbf{Y} \oslash \hat{\mathbf{Y}})}{(\mathbf{U}_N)^T} \quad (10)$$

259 where  $\odot$  is the element-wise multiplication,  $\oslash$  is the element-wise division and  
 260  $()^T$  is the transpose operator. The set of activation and basis matrices for both  
 261 the internal and external magnitude spectrograms is obtained updating the rules  
 262 detailed in Equations (6)-(10) using an iterative process until the algorithm  
 263 converges or reaches a maximum number of iterations  $M$ .

264 Focusing on the separation process applied to the biomedical sounds and  
 265 ambient noises present in the internal spectrogram, the estimated magnitude  
 266 spectrograms  $\hat{\mathbf{X}}_N$  and  $\hat{\mathbf{X}}_S$  can be obtained from the estimated basis  $\mathbf{U}_N, \mathbf{U}_S$   
 267 and activation matrices  $\mathbf{V}_N, \mathbf{V}_S$  as follows:

$$\hat{\mathbf{X}}_N = \mathbf{U}_N \mathbf{V}_N \quad (11)$$

268

$$\hat{\mathbf{X}}_S = \mathbf{U}_S \mathbf{V}_S \quad (12)$$

269 In order to denormalize the estimated magnitude spectrograms of the in-  
 270 ternal spectrogram, the matrices  $\hat{\mathbf{X}}_N, \hat{\mathbf{X}}_S$  are multiplied by the denominator of  
 271 Equation (1) when  $\mathbf{Z} = \mathbf{X}$ . To guarantee a conservative strategy in the re-  
 272 construction process, the estimated biomedical signal  $\hat{s}(t)$  (Equation (14)) is  
 273 computed by the inverse overlap-add STFT of the element-wise multiplication  
 274 between the complex spectrogram  $\mathbf{X}_c$  and a Wiener mask [11, 7] that represents  
 275 the relative energy contribution of the biomedical sounds to the energy of the  
 276 internal signal  $x(t)$ . The estimated ambient noise signal  $\hat{n}(t)$  (Equation (13)) is  
 277 calculated in a similar way as explained above in Equation (14), but now taking  
 278 into account that the Wiener mask explains the relative energy contribution of  
 279 the ambient noise sounds to the energy of the internal signal  $x(t)$ .

$$\hat{n}(t) = IDFT \left( \mathbf{X}_c \odot \frac{|\hat{\mathbf{X}}_N|^2}{\left( |\hat{\mathbf{X}}_N|^2 + |\hat{\mathbf{X}}_S|^2 \right)} \right) \quad (13)$$

280

$$\hat{s}(t) = IDFT \left( \mathbf{X}_c \odot \frac{|\hat{\mathbf{X}}_S|^2}{\left( |\hat{\mathbf{X}}_S|^2 + |\hat{\mathbf{X}}_N|^2 \right)} \right) \quad (14)$$

281 The pseudo code of the proposed method 2C-NMPCF for the ambient denois-  
 282 ing in auscultation is summarized in the Algorithm 1. Although the proposed  
 283 method can return the estimated biomedical signal  $\hat{s}(t)$  and the estimated am-  
 284 bient noise signal  $\hat{n}(t)$ , only the signal  $\hat{s}(t)$  is required for evaluation purposes  
 285 in this work.

---

**Algorithm 1** Ambient denoising using 2C-NMPCF

---

**Require:**  $y(t)$ ,  $x(t)$ ,  $K_N$ ,  $K_S$ ,  $\lambda$  and  $M$ .

- 1: Compute the normalized magnitude spectrogram  $\mathbf{X}$  of the internal signal  $x(t)$  using Equation (1).
  - 2: Compute the normalized magnitude spectrogram  $\mathbf{Y}$  of the external signal  $y(t)$  using Equation (1).
  - 3: Initialize each activation and basis matrix  $\mathbf{U}_N, \mathbf{U}_S, \mathbf{V}_N, \mathbf{V}_S, \mathbf{H}_N$  with random non-negative values.
  - 4: Update each activation and basis matrix  $\mathbf{U}_N, \mathbf{U}_S, \mathbf{V}_N, \mathbf{V}_S, \mathbf{H}_N$  using Equations (6)-(10) for the predefined number of iterations  $M$ .
  - 5: Compute the estimated magnitude spectrograms  $\hat{\mathbf{X}}_N$  using Equation (11).
  - 6: Compute the estimated magnitude spectrograms  $\hat{\mathbf{X}}_S$  using Equation (12).
  - 7: Denormalize the estimated magnitude spectrograms  $\hat{\mathbf{X}}_N, \hat{\mathbf{X}}_S$  multiplying by a factor equal to the denominator of Equation (1) when  $\mathbf{Z} = \mathbf{X}$ .
  - 8: Synthesize the estimated ambient noise  $\hat{n}(t)$  using the Equation (13).
  - 9: Synthesize the estimated biomedical signal  $\hat{s}(t)$  using the Equation (14).
- return**  $\hat{s}(t)$
- 

286 *2.3.4. Improving the sound quality of biomedical signals by means of an incre-*  
 287 *mental algorithm based on 2C-NMPCF*

288 The main limitation of the proposal 2C-NMPCF is related to the objective  
 289 function (see Equation (4)) used to minimize the residuals of the ambient noise.  
 290 Specifically, the iterative process based on the multiplicative update rules that

291 obtains the denoised biomedical basis and activation matrices is applied until  
 292 the convergence of 2C-NMPCF into a local minimum after  $M$  iterations. For  
 293 this reason, 2C-NMPCF by itself is not able to extract all spectral patterns  
 294 associated to ambient noises. To overcome the previous limitation, we propose  
 295 an incremental algorithm that runs 2C-NMPCF more than once improving the  
 296 estimated biomedical signal  $\hat{s}_i(t)$  obtained in the incremental iteration  $i$  by  
 297 removing additional spectral content associated to ambient noise that 2C-NMPCF  
 298 was not able to remove in the previous incremental iteration  $i - 1$ . In general  
 299 considering the iteration  $i$ , the internal signal  $x_{i+1}(t)$  of the next incremental  
 300 iteration  $i + 1$  is the estimated biomedical signal  $\hat{s}_i(t)$  obtained in the current  
 301 incremental iteration  $i$ , that is,  $x_{i+1}(t) = \hat{s}_i(t)$ . Note that in the first incremen-  
 302 tal iteration  $i = 1$ ,  $x_1(t) = x(t)$ . However, the external signal  $y(t)$  is fixed for all  
 303 incremental iterations since we assume that the signal  $y(t)$  is only composed by  
 304 ambient noises (no biomedical sounds are active). The flowchart of the proposed  
 305 incremental algorithm based on 2C-NMPCF is shown in Figure 4.

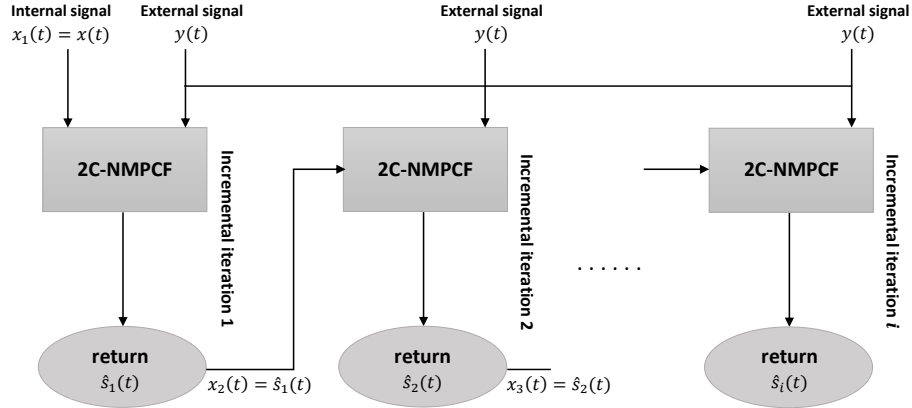


Figure 4: Flowchart of the proposed incremental algorithm based on 2C-NMPCF.

306 We assume the following assumptions in order to justify our incremental  
 307 proposal based on 2C-NMPCF: (i) the objective function  $\Gamma$  would converge into  
 308 a better local minimum at each incremental iteration since it would find re-  
 309 mainder spectral patterns of ambient noise that have not been extracted in the



310 previous iteration  $i - 1$  but they are still being repeated in both the internal  $\mathbf{X}_i$   
 311 and the external  $\mathbf{Y}$  magnitude spectrograms in the current incremental itera-  
 312 tion  $i$ ; (ii) 2C-NMPCF will remove most of the ambient noise while preserving  
 313 the content of the biomedical sounds until an optimal number of incremental  
 314 iterations  $i = i_o$ . From this optimal iteration  $i = i_o$ , our proposal can continue  
 315 to eliminate hidden patterns of ambient noise that are still active at the expense  
 316 of eliminating also spectral content related to biomedical signals. Summariz-  
 317 ing, this incremental approach attempts to maintain most of the biomedical  
 318 content  $\hat{s}_{i_o}(t)$  removing most of the ambient noise through the incremental iter-  
 319 ations. An illustrative example of the performance of the proposed incremental  
 320 approach is shown in Figure 5.

### 321 **3. Evaluation**

#### 322 *3.1. Metrics*

323 To evaluate the quality of the biomedical signals estimated by the proposed  
 324 method, the BSS EVAL toolbox [53, 54] has been used because it proposes a  
 325 set of metrics, widely applied in the field of sound source separation [11, 7] and  
 326 background noise removal [55], that quantify the quality of the sound separa-  
 327 tion between the original biomedical signal and its estimation. Two metrics,  
 328 measured in dB, are used as occurs in [56, 57]: (i) Source to Distortion ratio  
 329 (SDR) measures the overall quality of the estimated biomedical signal; and (ii)  
 330 Source to Interference ratio (SIR) measures the presence of ambient noise in the  
 331 estimated biomedical signal. Higher values of these ratios indicate better sound  
 332 quality of the estimated biomedical signal.

333 In this paper, the optimization and testing results have been obtained cal-  
 334 culating both SDR and SIR median values [58]. These results do not show  
 335 the absolute SDR and SIR values obtained from the estimated biomedical sig-  
 336 nal  $\hat{s}_i(t)$  but the SDR and SIR improvement comparing  $\hat{s}_i(t)$  and the original  
 337 internal signal  $x(t)$ .

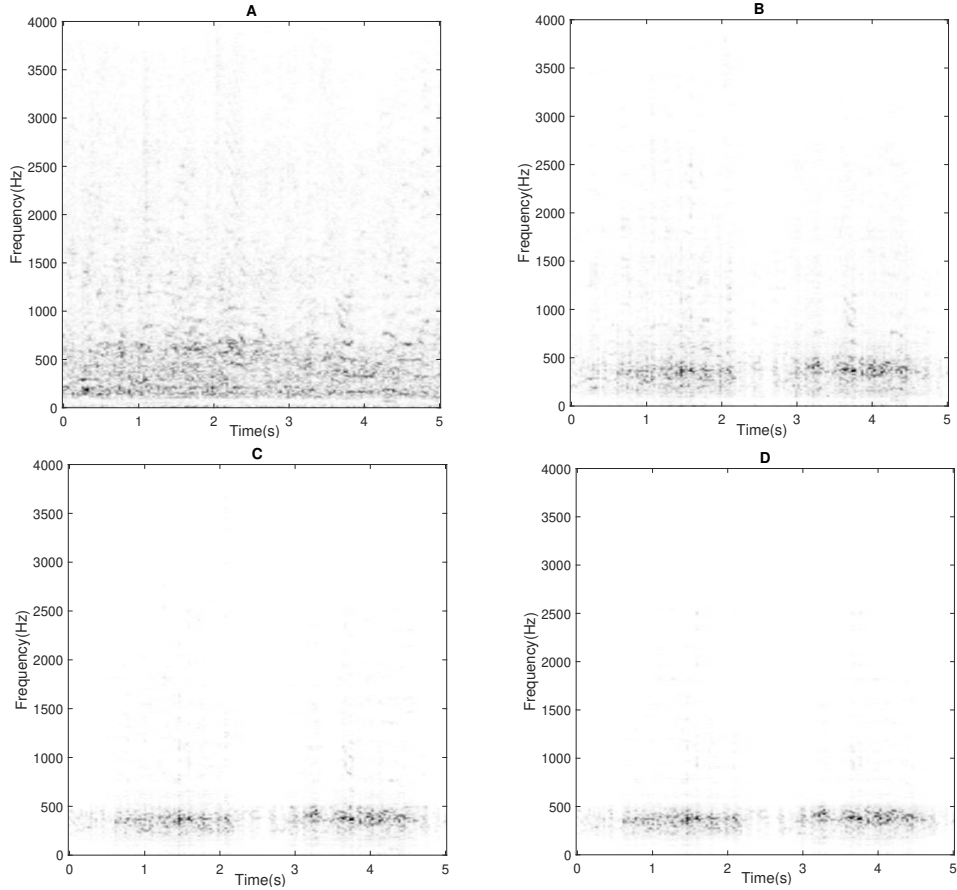


Figure 5: **A)** Magnitude spectrogram  $\mathbf{X}$  of the internal signal previously shown in Fig. 1 (Right); Some of the estimated biomedical magnitude spectrograms  $\hat{\mathbf{X}}_S$  provided by the incremental algorithm based on 2C-NMPCF through the incremental iterations  $i$ : **B)**  $i = 1$ , **C)**  $i = 2$  and **D)**  $i = 3$ . Here, it can be observed that the estimated biomedical spectrogram is refined after each incremental iteration by removing spurious ambient noise that are still active in the previous incremental iteration maintaining most of the biomedical spectral content.

338 As occurs in [24], two metrics related to speech intelligibility (SI) are added  
 339 to the objective assessment of the proposed method, which have been previously  
 340 used in the field of ambient noise suppression in lung auscultation [24, 25]: (i)  
 341 Normalized-Covariance Measure (NCM); and (ii) Coherence Speech Intelligi-  
 342 bility Index (CSII). Specifically, a three-level CSII approach is used by divid-  
 343 ing the signal into three amplitude regions: low ( $\text{CSII}_{low}$ ), mid ( $\text{CSII}_{mid}$ ) and

344 high (CSII<sub>high</sub>). In this work, each metric was computed between the origi-  
345 nal biomedical signal  $s(t)$  and the estimated biomedical signal  $\hat{s}(t)$ . Note that  
346 higher values of these metrics indicate better sound quality of the estimated  
347 biomedical signal. Finally, more details of these metrics can be found by the  
348 reader in [24, 25, 59].

### 349 3.2. Setup

350 Because most of the spectral content both the biomedical signals [45, 46, 48,  
351 49] and the ambient noise [23, 60] is concentrated below 4 KHz, in this work, a  
352 sampling rate equals  $f_s = 8$  KHz has been used as occurs in [24].

353 A preliminary study showed that the following parameters for time-frequency  
354 representation provide the best trade-off between the separation performance  
355 and the computational cost: a Hamming window with  $N = 512$  samples length  
356 (64ms) and 50% overlap; and a discrete Fourier transform using  $2N$  points  
357 similarly as in [11, 52]. Furthermore, the convergence of the proposed method  
358 was empirically observed after 50 iterations, so the parameter  $M$  is fixed to  
359  $M = 50$ .

360 Finally, note that the performance of the proposed method depends on the  
361 initial values with which the basis matrices  $\mathbf{U}_S$ ,  $\mathbf{U}_N$  and the activation matrices  
362  $\mathbf{V}_S$ ,  $\mathbf{V}_N$ ,  $\mathbf{H}_N$  have been initialised. Although the obtained results are not dis-  
363 persed and keep the same behavior, in order to overcome this issue, we have run  
364 the proposed method three times for each mixture and the results shown in this  
365 paper have been calculated using the median values as previously mentioned.

### 366 3.3. Results

367 In this section, experimental results related to the optimization and testing  
368 are detailed.

#### 369 3.3.1. Optimization results

370 Several parameters must be fitted to obtain the best performance of the  
371 proposed method in the removal of ambient noise. Four parameters have been

372 evaluated using the database  $D_O$ : (i) The number of biomedical bases  $K_S$   
373 used to characterise the spectral content of the biomedical signal  $s(t)$ , specif-  
374 ically,  $K_S \in [16, 32, 64, 128, 256]$ ; (ii) The number of ambient noise bases  $K_N$   
375 used to characterise the spectral content of the ambient noise  $n(t)$ , specifically,  
376  $K_N \in [16, 32, 64, 128, 256]$ ; (iii) The value  $\lambda$  to balance the importance of the in-  
377 ternal  $\mathbf{X}$  and external  $\mathbf{Y}$  magnitude spectrograms in the co-factorization process.  
378 In this case,  $\lambda \in [0.01, 0.1, 1, 10, 25, 50, 100, 250, 500, 1000]$ ; (iv) The number of  
379 incremental iterations  $i$  applied to 2C-NMPCF.

380 The optimization process is composed of two steps:

- 381 1. Step I. Optimize the three parameters  $K_S, K_N, \lambda$  in order to reach the  
382 greater median of the SDR improvement when applying 2C-NMPCF con-  
383 sidering all the types of ambient noises and SNR previously mentioned.
- 384 2. Step II. Once the parameters of 2C-NMPCF have been optimized, it must  
385 be obtained the optimal number of incremental iterations  $i = i_o$  to achieve  
386 the best performance of the proposed method (see Figure 4).

387 Figure 6 shows the median of the SDR improvement analyzing the search  
388 space derived from the parameters  $\lambda, K_S$  and  $K_N$ . Results indicate that the  
389 proposed method provides the best denoising performance, by means of the  
390 maximum median value of the SDR improvement, using the optimal param-  
391 eters  $K_N=256$  and  $K_S=16$ . These optimal values demonstrate that ambient  
392 noise requires a greater number of spectral patterns compared to biomedical  
393 sounds due to their greater spectral diversity. The analysis of different lung and  
394 heart signals indicates that the spectral modeling of these biomedical sounds  
395 is simpler and therefore needs a smaller dictionary of bases since lung sounds  
396 could be factorize by a low-rank decomposition using broadband spectral pat-  
397 terns that show temporal and spectral smoothness. However, heart sounds could  
398 be modeled as low-frequency pulses located in regular intervals in time.

399 Figure 7 shows the median of the SDR and SIR improvement using the  
400 previous optimal values of the parameters  $K_S$  and  $K_N$  (that is,  $K_S=16$  and  
401  $K_N=256$ ). It can be confirmed that giving importance to the sharing of spectral

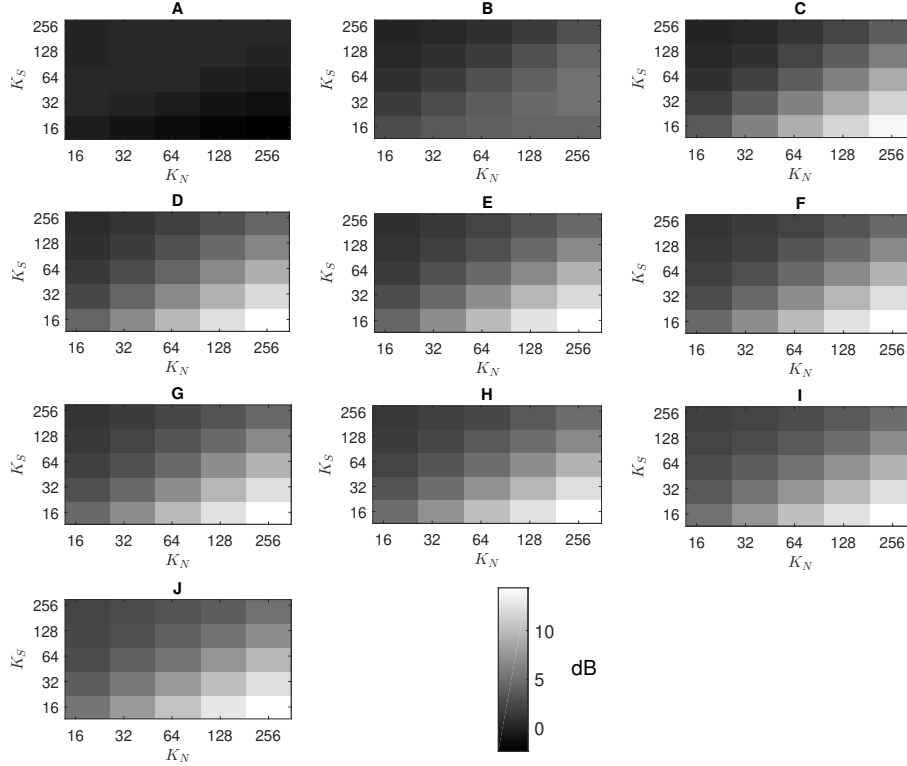


Figure 6: Median values of the SDR improvement (dB) evaluating  $D_O$  for the following  $\lambda$  values:  $\lambda=0.01$  (A),  $\lambda=0.1$  (B),  $\lambda=1$  (C),  $\lambda=10$  (D),  $\lambda=25$  (E),  $\lambda=50$  (F),  $\lambda=100$  (G),  $\lambda=250$  (H),  $\lambda=500$  (I) and  $\lambda=1000$  (J).

402 bases in the joint factorization process finds a better local minimum in the fac-  
 403 torization process, since ambient noise clearly reveals its simultaneous presence  
 404 both in the internal and external signal spectrogram. Results report a signif-  
 405 icant and stable SDR and SIR improvement equals 14 dB and 19.5 dB using  
 406  $\lambda \geq 10$ . For this reason, we have chosen the optimal parameter  $\lambda=10$ .

407 Figure 8 depicts the optimal number of incremental iterations  $i$  of the pro-  
 408 posed method using the previous optimal parameters  $K_S$ ,  $K_N$  and  $\lambda$ . It is  
 409 observed that both SDR and SIR improvement increases sharply when apply-  
 410 ing 2C-NMPCF in the second incremental iteration ( $i=2$ ), higher increase for  
 411 SIR compared to SDR. In the third incremental iteration ( $i=3$ ), the SDR im-

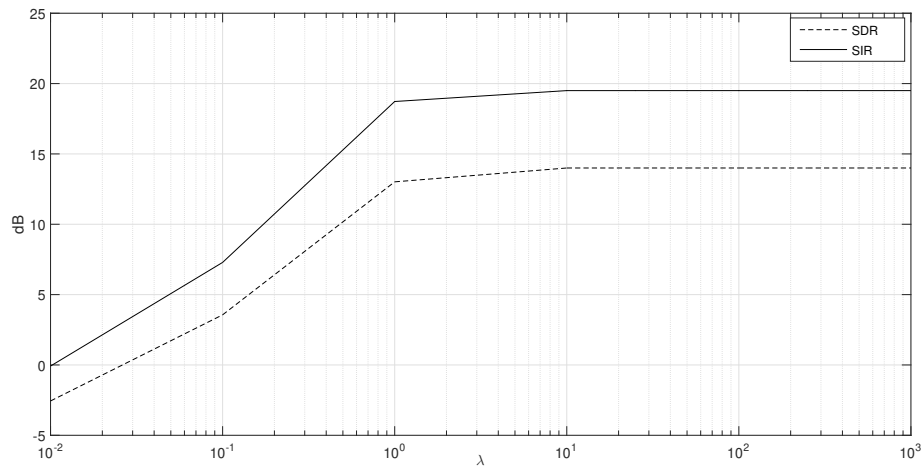


Figure 7: Median values of the SDR (dashed line) and SIR (solid line) improvement of the proposal algorithm evaluating  $D_O$ , keeping fixed  $K_S=16$  and  $K_N=256$  varying the parameter  $\lambda$ .

412 provement increases slightly and then starts to decrease gradually while the  
 413 SIR improvement continues to grow. Experimental results indicate that am-  
 414 bient noise continues to be suppressed at the expense of starting to remove  
 415 biomedical spectral content when  $i > 3$ . For this reason, the optimal number  
 416 of incremental iterations has been set at  $i_o = 3$  with the aim of providing the  
 417 greatest suppression of ambient noise while maintaining most of the biomed-  
 418 ical content, being  $i_o$  the iteration in which the maximum SDR improvement is  
 419 obtained.

### 420 3.3.2. Objective results simulating an ideal scenario

421 This section evaluates the ambient denoising performance of the proposed  
 422 method simulating an ideal scenario since neither the effects of propagation on  
 423 the patient’s body material nor the acoustics of the room are considered active  
 424 (these effects will be analyzed in section 3.3.3.).

425 Figure 9 shows SDR and SIR improvement comparing the behaviour of the  
 426 ambient noise removal evaluating the database  $D_T$  for the proposed method  
 427 (2C-NMPCF) and the baseline methods MSS and NLMS. Hereinafter, the SDR

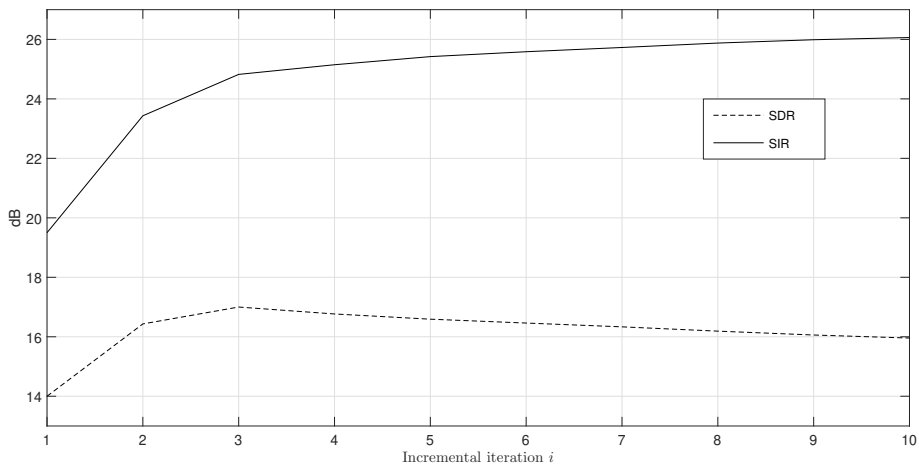


Figure 8: Median values of the SDR (dashed line) and SIR (solid line) improvement with the number of incremental iterations  $i$  evaluating  $D_O$ , keeping fixed  $(K_S, K_N, \lambda)=(16, 256, 10)$ .

428 and SIR improvement associated to each method are represented by  $SDR_P$  and  
 429  $SIR_P$  (the proposed method),  $SDR_M$  and  $SIR_M$  (MSS) and finally,  $SDR_N$  and  
 430  $SIR_N$  (NLMS). Each box represents 250 data points, one for each recording of  
 431 the database  $D_T$ . The lower and upper lines of each box show the 25th and  
 432 75th percentiles. The line in the middle of each box represents the median value.  
 433 The diamond in the center of each box represents the average value. The lines  
 434 extending above and below each box show the extent of the rest of the samples,  
 435 excluding outliers. Outliers are defined as points that are over 1.5 times the  
 436 interquartile range from the sample median, which are shown as crosses.

437 Figure 9A shows that the proposed method outperforms MSS and NLMS, in  
 438 terms of SDR, in most of the SNR, obtaining the best average performance when  
 439 all SNRs are taken into account. It can be seen how the proposed method reports  
 440 the most robust SDR performance compared to MSS and NLMS, mainly in high  
 441 noise environments ( $SNR \in [-20 \text{ dB}, -10 \text{ dB}]$ ), by means of a stable SDR trend in  
 442 this SNR range. It suggests that the multichannel co-factorization approach is  
 443 less dependent on the ratio of ambient noise to biomedical content than the other  
 444 MSS and NLMS methods. Although MSS shows an increasing trend in SDR  
 445 improvement stabilizing its denoising results for  $SNR \geq -10 \text{ dB}$ , these results are

446 not satisfactory enough to consider MSS as a competitive method with respect  
447 to the proposed method or NLMS since the same does not happen in terms of  
448 SIR improvement (see Figure 9B) where MSS provides  $SIR_M \leq 14\text{dB}$  in high  
449 noise environments compared mainly to the proposed method.

450 Figure 9B indicates that although the proposed method and NLMS obtain a  
451 very similar performance considering the average of all SNRs, NLMS is slightly  
452 better in  $SNR \leq -10\text{ dB}$  at the expense of losing a large amount of biomedical  
453 content, a fact that does not occur with the proposed method. This behavior  
454 shown by NLMS is confirmed by its SDR and SIR reduction as SNR increases.  
455 Focusing on the proposed method, the stable SDR and SIR trends regardless  
456 of the SNR evaluated demonstrate that our incremental approach is a more  
457 reliable feature to remove ambient noise because: i) MSS provides a satisfactory  
458 performance assuming a distortion of biomedical sounds at low frequencies and  
459 penalizing the occurrences of noise with strong energy in high spectral bands  
460 [24]; and ii) the proposed method assumes ambient noise as a repetitive sound  
461 event found in both the internal and external spectrograms, so it does not  
462 depend on any temporal misalignment between the internal and external signal  
463 as occurs with NLMS.



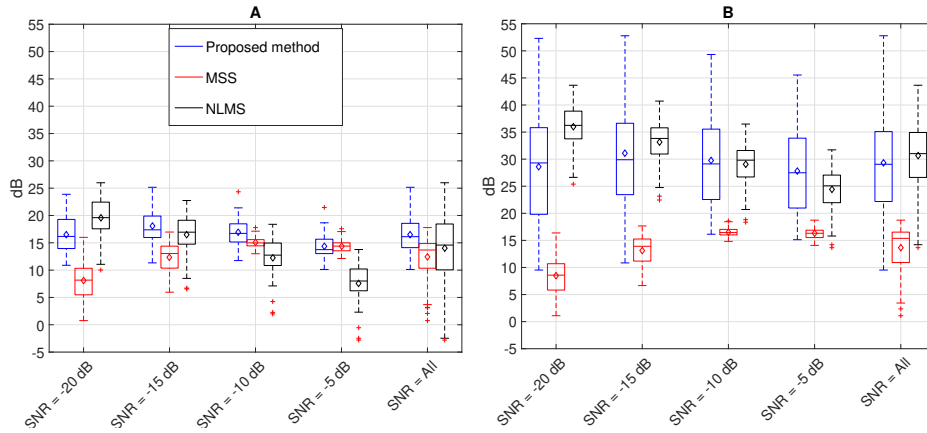


Figure 9: SDR (A) and SIR (B) improvement considering all the noises (Ambulance Siren, Baby Crying, Babble, Car and Street) for each SNR from database  $D_T$ . Each box represents the denoising performance results for each method evaluated. The color of the legend, associated to each evaluated method shown in the subfigure 9.A, refers to the same methods for all subfigures.

464 Figure 10 shows a detailed analysis of the denoising performance of the  
 465 proposed method considering each particular type of ambient noise evaluated  
 466 previously in Figure 9. Each box represents 50 data points, one for each recording  
 467 of the database  $D_T$ . Highlight the stable SDR and SIR trends regardless  
 468 of the type of noise and level SNR evaluated shown by the proposed method  
 469 unlike what happens with MSS and NLMS.

- 470 • Ambulance siren noise: Figures 10A and 10B show that although NLMS  
 471 only obtains better SDR results compared with the proposed method in  
 472 high noisy scenarios ( $\text{SNR} \leq -10$  dB), the denoising performance of both  
 473 methods are still competitive averaging all evaluated SNRs. Nevertheless,  
 474 it is interesting that the SDR improvement of the proposed method out-  
 475 performs NLMS in acoustic scenarios in which  $\text{SNR} \geq -10$  dB because it  
 476 is a frequent acoustic scenario that can be found around an auscultation  
 477 room located inside a health center.
- 478 • Baby crying noise: Figures 10C and 10D show that the proposed method  
 479 significantly improves the SIR denoising performance compared with MSS

480 and NLMS keeping most of the biomedical content as shown SDR results,  
481 an interesting advantage particularly in high noisy environments. Results  
482 seem to suggest that the strong harmonic structure contained into this  
483 type of ambient noise facilitates the sound separation between noise and  
484 biomedical sounds since the more dissimilar the spectral patterns of the  
485 noise and the biomedical signal, the better the noise suppression perfor-  
486 mance of the proposed method.

- 487 • Babble noise: Figures 10E and 10F indicate that NLMS obtains the best  
488 SDR and SIR improvements evaluating this particular ambient noise but  
489 showing a decreasing trend as the SNR increases, unlike the other eval-  
490 uated methods. Experimental results indicate that the spectrum of the  
491 Babble noise and the biomedical signal, mainly lung sounds, is more sim-  
492 ilar compared to the above ambient noises (ambulance siren and baby  
493 crying noise) so, the greater are the spectral differences between the tar-  
494 get sounds and the noise, the better is the performance of the proposed  
495 method since the repetitive behaviour of the ambient noise in the co-  
496 factorization process is more easy to suppress.
- 497 • Car noise: Figures 10G and 10H indicate that the proposed method  
498 obtains the best ambient denoising performance compared to MSS and  
499 NLMS since it achieves the best biomedical audio quality by maximizing  
500 the amount of suppressed ambient noise at the expense of preserving most  
501 of the biomedical content related to the estimated biomedical signal of  
502 interest.
- 503 • Street noise: Figures 10I and 10J show that NLMS preserves higher  
504 amount of biomedical content in high noisy environments ( $\text{SNR} \leq -10$  dB)  
505 compared to the proposed method and MSS. However, SDR and SIR re-  
506 sults indicate that the proposed method and MSS are competitive com-  
507 pared to NLMS. Moreover, NLMS reports the first SIR ranking position  
508 in all SNR scenarios. Finally, the proposed method outperforms, in terms  
509 of SDR results, MSS and NLMS in the remainder SNRs.

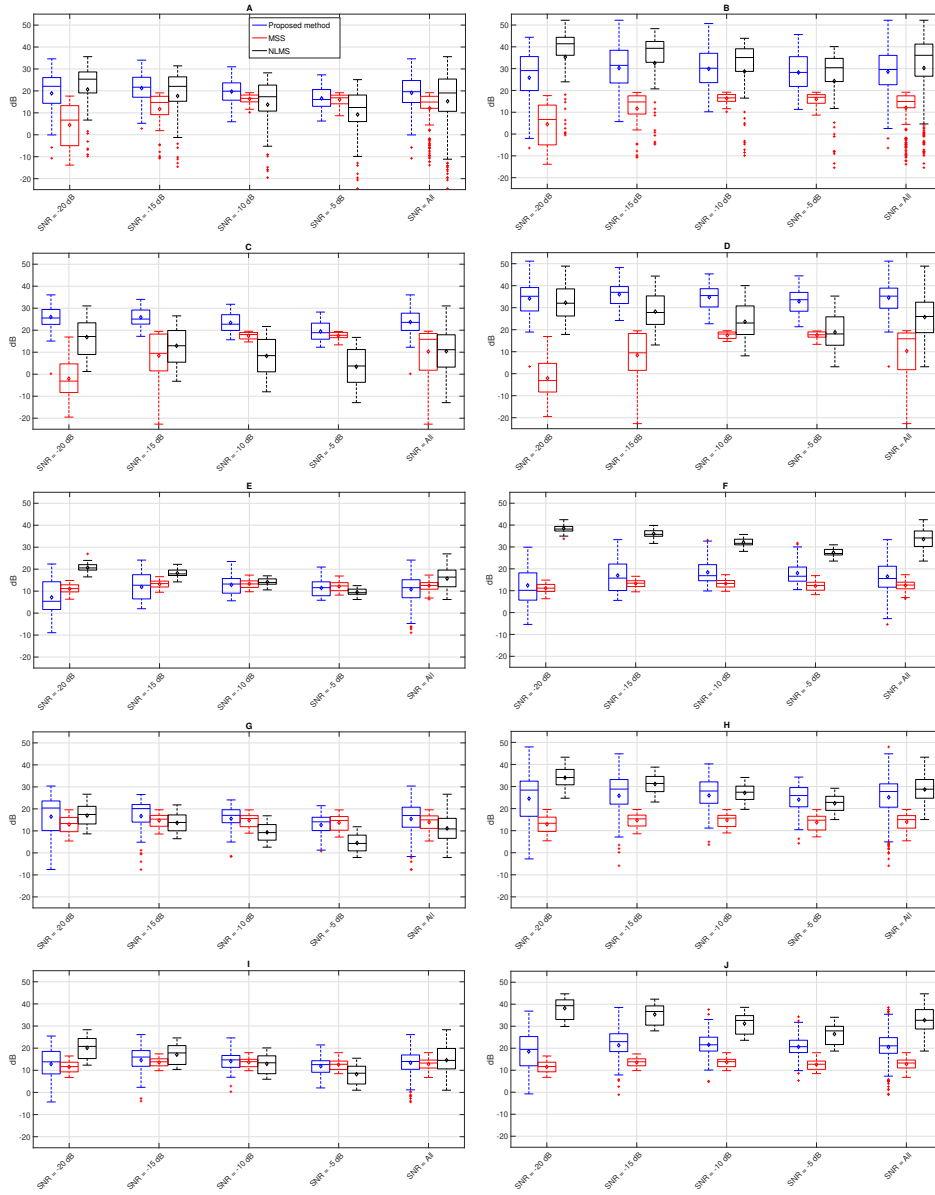


Figure 10: SDR and SIR improvement results provide by the proposed method and the baseline methods (MSS and NLMS) evaluating  $D_T$  and each type of ambient noise along SNR: Ambulance Siren (A and B), Baby crying (C and D), Babble (E and F), Car (G and H) and Street (I and J). Each box represents the denoising performance results for each method evaluated. The color of the legend, associated to each evaluated method shown in the subfigure 10.A, refers to the same methods for all subfigures.

510 Figure 11 shows the effect of inserting a time delay between the internal and  
 511 external signals. The purpose of the delay is to simulate the time processing  
 512 that takes the digital stethoscope to apply filtering, artifacts removal and other  
 513 signal processing operations [61]. It can be observed how the delay affects  
 514 to the computation of the median value of the SDR and SIR improvement  
 515 evaluating all the previous ambient noises and SNRs. Results confirm that the  
 516 most remarkable advantage of the proposed method is its robustness with the  
 517 variation of the delay. In fact, the proposed method shows a stable behavior in  
 518 relation to the delay variation between the internal and external signals used  
 519 in the co-factorization, in contrast to a higher dependence of MSS compared to  
 520 the proposed method and a higher dependence of NLMS compared to MSS. The  
 521 ambient denoising performance between the proposed method, MSS and NLMS  
 522 is accentuated when the delay is active since SDR and SIR results obtained  
 523 by MSS and NLMS are reduced as the delay increases. Comparing with MSS  
 524 and NLMS, the proposed method obtains an improvement of 13.1 dB and 16.5  
 525 dB in terms of SDR, and 21 dB and 21.5 dB in terms of SIR applying a delay  
 526 equals to 25 milliseconds. As previously mentioned, results confirm the proposed  
 527 method as a more appropriate approach to remove the ambient noise because  
 528 the multichannel co-factorization is based of the noise modelling as repetitive  
 529 spectral patterns that can be found in any time of the internal and external  
 530 spectrograms, avoiding errors due to presence of temporal misalignment between  
 531 both spectrograms.

532 Figure 12 shows the metrics NCM,  $CSII_{low}$ ,  $CSII_{med}$  and  $CSII_{high}$  results  
 533 comparing the performance of the ambient noise removal evaluating the database  
 534  $D_T$  for the proposed method, MSS and NLMS. The higher metric value is ob-  
 535 tained, greater acoustic similarity between the estimated biomedical signal and  
 536 the original biomedical signal is with respect to their sound contents as occurs  
 537 in [24]. Specifically, each box represents 250 data points, one for each record-  
 538 ing of the database  $D_T$ . Figure 12 reports that, in general, the best average  
 539 ambient denoising results are obtained by NLMS. However, it is observed that  
 540 the average ambient noise suppression achieved by the proposed method can be

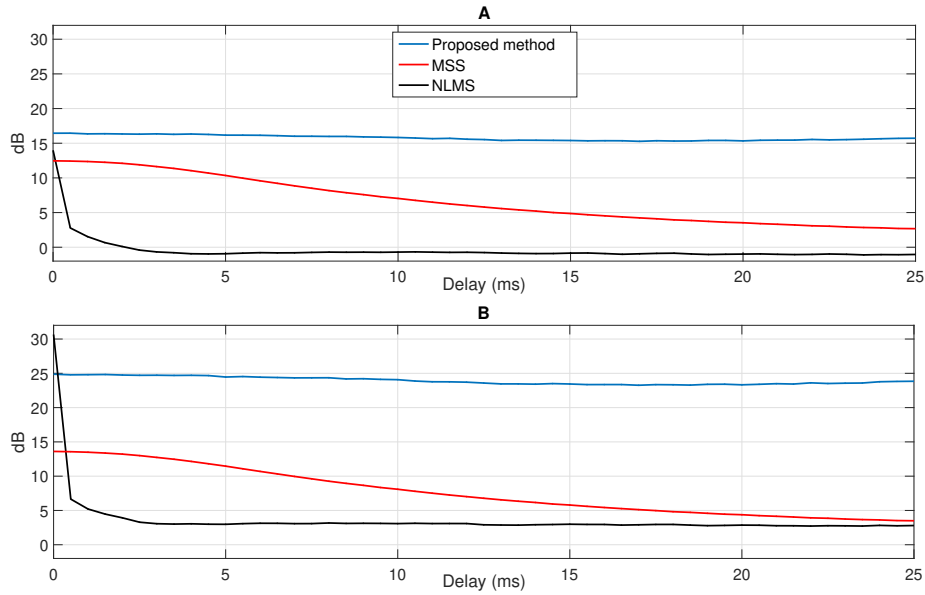


Figure 11: Median values of the SDR (A) and SIR (B) improvement analyzing all the noises (Ambulance Siren, Baby Crying, Babble, Car and Street) and SNRs (-20 dB, -15 dB, -10 dB, -5 dB) varying the delay between the internal and external spectrograms in the database  $D_T$ . The color of the legend, associated to each evaluated method shown in the subfigure 11.A, refers to the same methods for all subfigures.

541 considered very similar to MSS in most cases of SNR. Broadly, results obtained  
 542 by each evaluated method improve the acoustic quality of biomedical sounds by  
 543 eliminating ambient noise that typically hinders clinical examination. Finally,  
 544 it can be seen how the above results fall as the SNR decreases, similar to what  
 545 happens in the real world because it is more difficult to hear biomedical sounds  
 546 in those acoustic scenarios where biomedical sounds are barely audible due to  
 547 high ambient noise levels.

### 548 3.3.3. Objective results simulating a real scenario

549 This section assesses the ambient denoising performance of the proposed  
 550 method simulating a real scenario where both the propagation in the patient's  
 551 body material and the acoustics of the room have been considered active.

552 From the noise recordings belonging to the database  $D_N$  (see section 2.1),

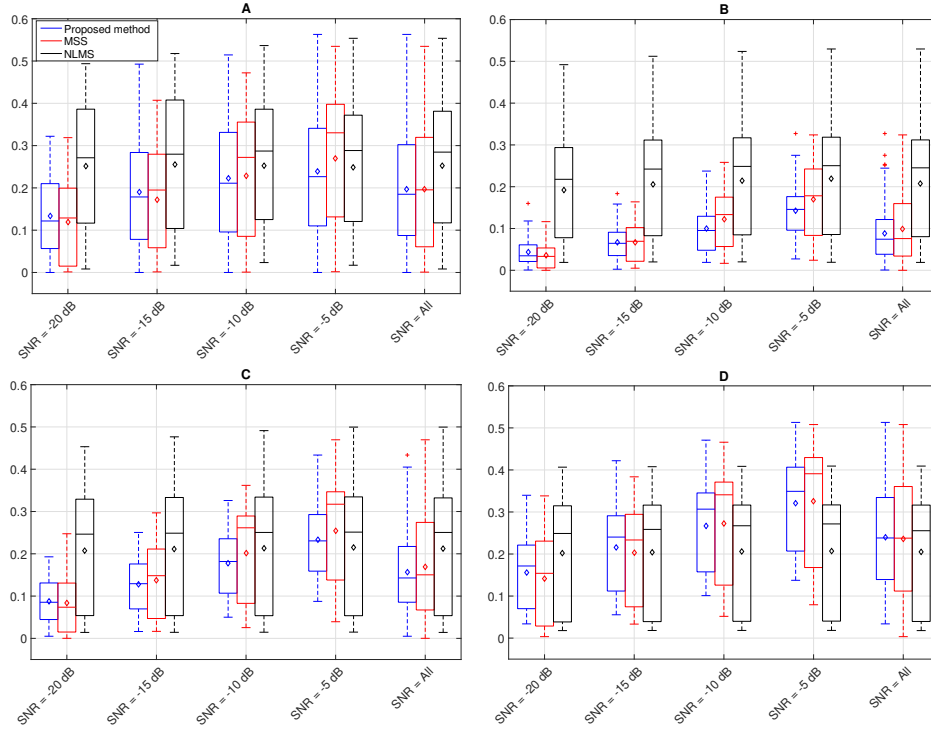


Figure 12: NCM (A),  $CSII_{low}$  (B),  $CSII_{med}$  (C) and  $CSII_{high}$  (D) results considering all the noises (Ambulance Siren, Baby Crying, Babble, Car and Street) for each SNR from database  $D_T$ . Each box represents the denoising performance results for each method evaluated. The color of the legend, associated to each evaluated method shown in the subfigure 12.A, refers to the same methods for all subfigures.

553 we have selected those indoor and outdoor ambient noises that can be heard  
 554 surround a medical consultation room located in a health center: ambulance  
 555 siren [35, 36], baby crying [37], babble (people speaking) [38, 39] and street (car  
 556 passing by, car engine running, car idling, bus, truck, children yelling, people  
 557 talking, workers on the street) [41, 42], that is, a total of 120 single-channel  
 558 recordings of noises (30 recordings per type of noise) lasting 5 seconds each  
 559 noise recording as occurred in section 2.1.

560 The database  $D_F$  has been created by modifying the database  $D_T$  in order to  
 561 evaluate the ambient denoising performance of the proposed method taking into  
 562 account the effects of propagation on the materials of the patient's body and the

563 acoustics of the room. Specifically, two impulse responses, from the open-access  
564 dataset [62, 63] and measured with an adult human, have been applied to the  
565 previous database  $D_T$ . The first impulse response  $h_H(t)$  has been calculated  
566 locating the microphone in the position 33 [62] because it can be considered a  
567 correct placement to replicate a heart auscultation in real conditions. The sec-  
568 ond impulse response  $h_L(t)$  has been calculated locating the microphone in the  
569 position 55 [62] because it can be considered a correct placement to replicate a  
570 lung auscultation in real conditions. Both impulse responses have been resam-  
571 pled at 8kHz and calculated using a sound source located behind the chest of the  
572 human. More details related to the propagation on the patient's body material  
573 can be found in [62]. In order to simulate the acoustics of the room, we have  
574 assumed a standard room in the Hospital Universitario de Jaen (Spain) with  
575 dimensions of 7m of large, 4m of width and 2.7m of height using the well-known  
576 image method [64] and the MATLAB RIR generator<sup>1</sup>. Specifically, the room  
577 impulse response  $h_R(t)$  has been designed assuming a moderate reverberation  
578 time  $RT_{60}$  of 0.4 seconds. The sensor is an omnidirectional microphone placed  
579 at the center of the room and the sources are placed outside the room, one in  
580 the waiting room and the other outdoors. As occurs with the database  $D_T$ ,  
581 several SNR have been applied in the mixing process to create the database  
582  $D_F$  in order to evaluate high noisy environments. In this way, the databases  
583  $D_{F_{-20}}$  (SNR=-20 dB),  $D_{F_{-15}}$  (SNR=-15 dB),  $D_{F_{-10}}$  (SNR=-10 dB) and  $D_{F_{-5}}$   
584 (SNR=-5 dB) refer to the same database  $D_F$  but using different SNR between  
585 biomedical and ambient noise recordings by means of the mixing process shown  
586 in Figure 13. The mixing process is performed as follows: i) each ambient noise  
587 signal  $n_i(t)$  (from the database  $D_N$ ) is convoluted with the impulse response  
588  $h_R(t)$  to create the external signal  $y(t)$ . Here, the signal  $y(t)$  is only affected  
589 by the acoustics of the room; ii) the convolution of the external signal  $y(t)$  and  
590 the impulse response of the human body,  $h_H(t)$  or  $h_L(t)$ , generates the ambi-  
591 ent noise signal  $n_{iRB}(t)$ . Here, the signal  $n_{iRB}(t)$  is affected by both the room

---

<sup>1</sup><https://www.audiolabs-erlangen.de/fau/professor/habets/software/rir-generator>

592 acoustics and the patient's body material; iii) the ambient noise signal captured  
 593 by the stethoscope  $n_{iMIX}(t)$  corresponds to the sum of  $n_{iRB}(t)$  and  $y(t)$ , that  
 594 is,  $n_{iMIX}(t)=n_{iRB}(t)+y(t)$ ; and iv) a desired SNR, in decibels, is obtained be-  
 595 tween the biomedical signal  $s(t)$  and the ambient noise signal  $n(t)$ . The signal  
 596  $n(t)$  is calculated weighting  $n_{iMIX}(t)$  by the parameter  $\alpha = 10^{\frac{(SNR_{init}-SNR)}{20}}$ ,  
 597 being  $SNR_{init}$  the initial signal-to-noise ratio between the biomedical signal  
 598  $s(t)$  and the signal  $n_{iMIX}(t)$ .

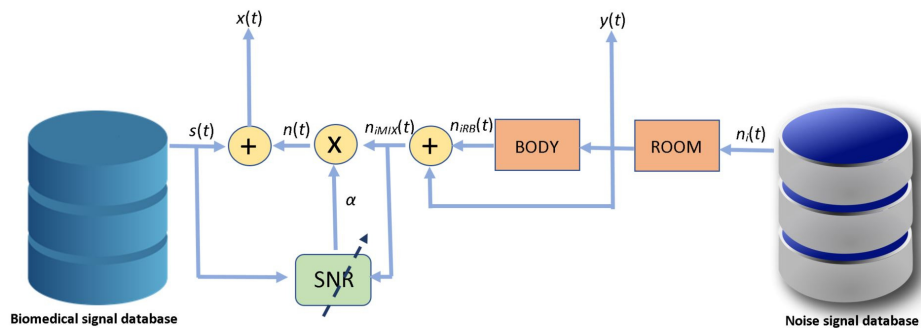


Figure 13: Scheme to simulate a more realistic scenario in which both the effects of propagation on the patient's body material and the acoustics of the room are active.

599 Figure 14 shows that the proposed method obtains the best average results  
 600 for SDR (16 dB additional with respect to MSS and 19 dB additional with re-  
 601 spect to NLMS) and SIR (23 dB additional with respect to MSS and 13 dB  
 602 additional with respect to NLMS), significantly outperforming the results pro-  
 603 vided by MSS and NLMS for each SNR evaluated. Once again, the stable trend  
 604 of SDR and SIR achieved by the proposed method for each SNR in compar-  
 605 ison with the rest of the methods is remarkable. The SDR results indicate  
 606 that neither MSS nor NLMS can be considered competitive methods in this  
 607 scenario since the results obtained by both methods yield mean values of SDR  
 608  $\leq 1$  dB, which implies a huge loss of biomedical content despite the fact that  
 609 NLMS gets a great isolated biomedical signal, showing mean values of  $SIR_N \in$   
 610 [8 dB, 14dB] but lower than those obtained by the proposed method for each  
 611 SNR. Moreover, the worst performance is provided by MSS in high noisy en-



612 vironments (SNR<-10dB) because the large amount of noise mixed with the  
 613 biomedical signal hinders optimal MSS performance by avoiding to estimate  
 614 a correct SNR-dependent factor. Nevertheless, these results again confirm the  
 615 better performance of NLMS compared to MSS similar to what occurs in Figure  
 616 9.

617 Comparing Figure 9 and Figure 14, it is evidenced the high robustness of  
 618 the proposed method compared to MSS and NLMS since the average reduction  
 619 related to the SDR and SIR improvement provided by the proposed method  
 620 taking into account the effects of this scenario is  $\leq 1$  dB and  $\leq 5$  dB while  
 621 the same SDR and SIR reduction increases to 14 dB and 15 dB for MSS, and  
 622 18 dB and 21 dB for NLMS. Figure 14 reports that the estimated denoised  
 623 spectrogram that models the ambient noise as a repetitive sound event in a  
 624 cofactorized NMF where simultaneously active repetitive spectral patterns are  
 625 sought in two different spectrograms provides higher sound quality compared to  
 626 the spectral subtraction approach or the adaptive filtering technique where the  
 627 use of temporal information between both spectrograms is essential to correctly  
 628 recover the target signal.

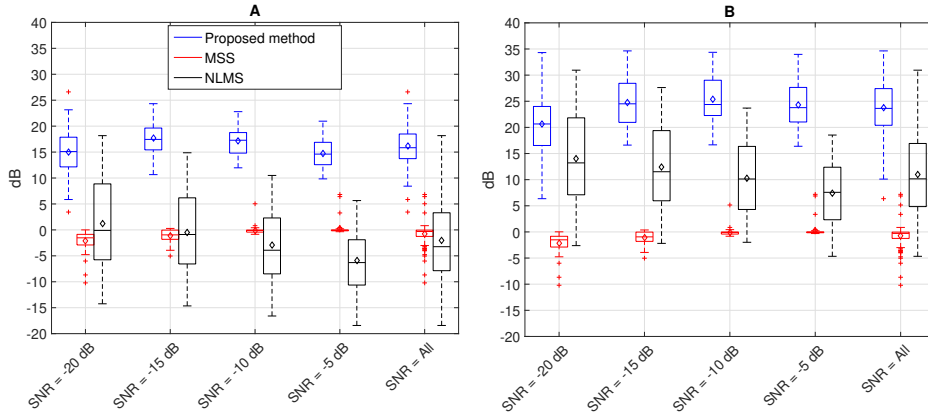


Figure 14: SDR (A) and SIR (B) improvement considering all the noises (Ambulance Siren, Baby Crying, Babble and Street) for each SNR from database  $D_F$ . Each box represents the denoising performance results for each method evaluated. The color of the legend, associated to each evaluated method shown in the subfigure 14.A, refers to the same methods for all subfigures.

629 Figure 15 evidences that the proposed method provides the best improve-  
630 ments of SDR and SIR for each type of ambient noise in most SNR by means  
631 of maximizing the denoising of ambient noise (highest  $SIR_P$ ) at the expense of  
632 minimizing the loss of biomedical energy (highest  $SDR_P$ ). It is also confirmed,  
633 as occurs in Figure 10, that the proposed method shows a stable trend in SDR  
634 and SIR unlike what happens with the other methods evaluated. Comparing  
635 Figure 15 and Figure 10, it can be deduced that the effect of sound propagation  
636 through the patient's body material and room acoustics reduces the denoising  
637 performance, in terms of SDR and SIR, in all the evaluated methods, however,  
638 this reduction is significantly lower in the proposed method compared to MSS  
639 and NLMS since the performance of MSS drops sharply in this acoustic environ-  
640 ment obtaining  $SDR_M$  and  $SIR_M$  results  $\leq 0$  dB, being the proposed method  
641 and NLMS (only for the case of street noise) the only competitive methods.  
642 Again, these results demonstrate what was previously stated in section 3.3.2,  
643 that the non-dependence of the temporal information, as occurs in the proposed  
644 method, between both spectrograms associated with the internal and external  
645 signal achieves a more adequate behavior in the denoising of environmental noise  
646 in comparison with those methods (MSS and NLMS) in which the temporal in-  
647 formation is used in greater or less way.

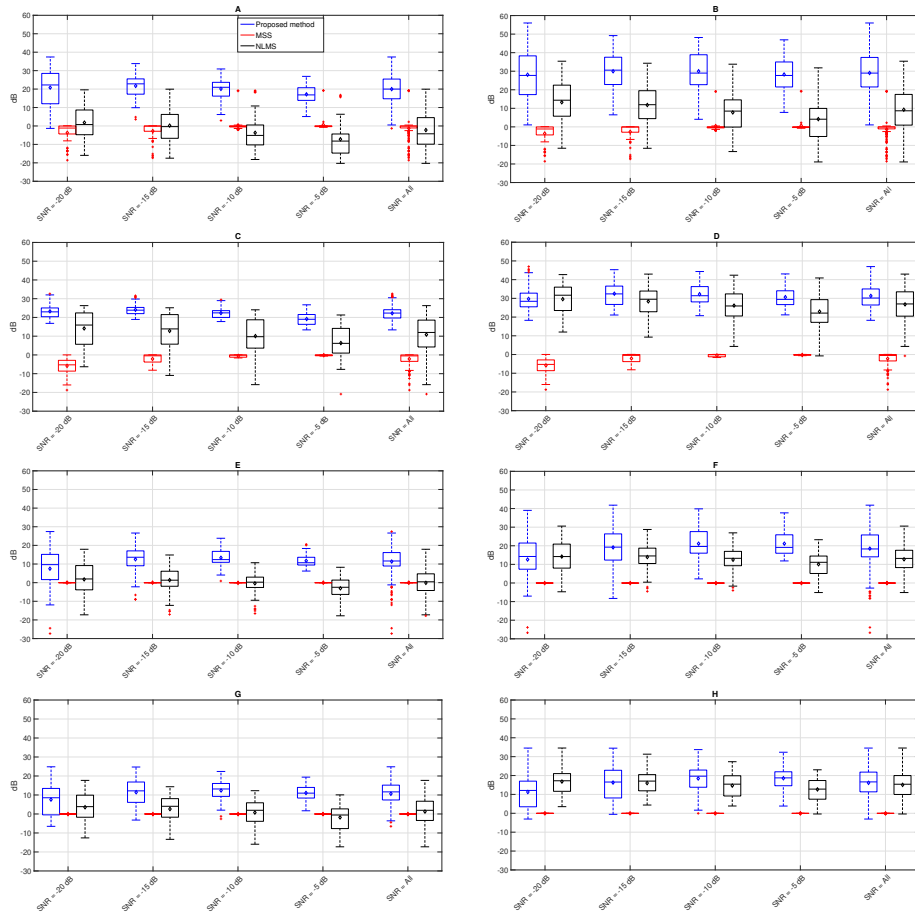


Figure 15: SDR and SIR improvement results provide by the proposed method and the baseline methods (MSS and NLMS) evaluating  $D_F$  and each type of ambient noise along SNR: Ambulance Siren (A and B), Baby crying (C and D), Babble (E and F) and Street (G and H). Each box represents the denoising performance results for each method evaluated. The color of the legend, associated to each evaluated method shown in the subfigure 15.A, refers to the same methods for all subfigures.

648        Figure 16 shows that the proposed method obtains the highest metrics NCM,  
649         $CSII_{low}$ ,  $CSII_{med}$  and  $CSII_{high}$  for all SNR evaluated increasing the denoising  
650        performance as the SNR increases. NLMS is ranked in second position followed  
651        by MSS in the last position. The ambient denoising performance provided by  
652        MSS (Figure 12) slightly exceeds the proposed method but this performance is

653 not the same when the propagation of the patient's body material and room  
654 acoustics are active (Figure 16). Comparing Figure 12 and Figure 16 and unlike  
655 what happens with the average values of NLMS, it is observed that the proposed  
656 method shows a clear increasing trend in these metrics as the SNR increases. As  
657 a result, it seems that the acoustic quality of the biomedical content estimated  
658 by the proposed method using these metrics is influenced to a greater extent  
659 by the ratio between the biomedical signal and the ambient noise in the input  
660 mixture, as opposed to its stable trend previously reported for SDR and SIR.  
661 Although the denoising performance provided by MSS and NLMS in Figure 16  
662 drops significantly for all evaluated metrics and SNR compared to the results  
663 obtained in Figure 12, the same results obtained by the proposed method are  
664 also reduced but this reduction can be considered marginal compared to that  
665 suffered by both MSS and NLMS. As a result, it can be deduced that the  
666 temporal non-dependence between the internal and external spectrograms in a  
667 multichannel co-factorization is a suitable characteristic to remove ambient noise  
668 in biomedical sounds mainly considering the effects typically found to model a  
669 real scenario (e.g., propagation of the patient's body material, reverberations  
670 in the room, ...). These effects drastically reduce the acoustic quality of the  
671 extracted biomedical signal as well as the robustness of the method since any  
672 temporal misalignment or energy changes caused by impulse responses related  
673 to the previous effects imply a relevant signal distortion in the spectrograms,  
674 being those approaches based on spectral subtraction (MSS) or adaptive filtering  
675 (NLMS) less efficient to minimize the sound interference caused by the previous  
676 ambient noises.

#### 677 *3.3.4. Computational complexity assessment*

678 In this section, the analysis of the computational complexity of each method  
679 evaluated considering the number of elementary operations (multiplications and  
680 additions) is detailed in Table 2. The number of operations per second can be  
681 obtained taking into account the parameters on which each algorithm used in the  
682 experimental results of this work is based: i) the proposed method with  $N=512$

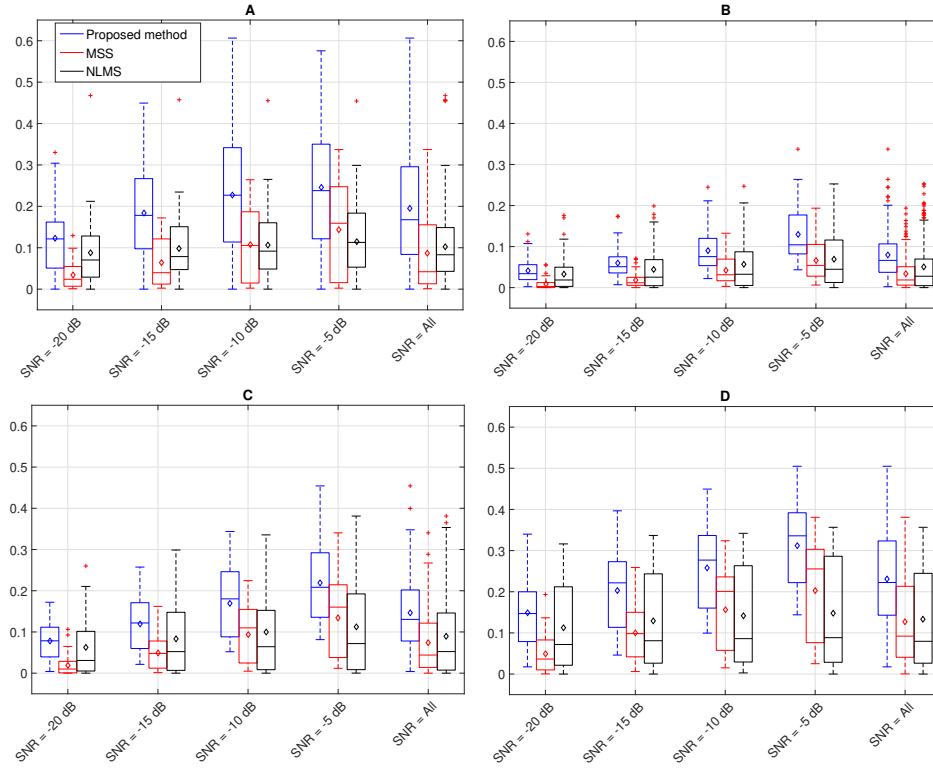


Figure 16: NCM (A),  $CSII_{low}$  (B),  $CSII_{med}$  (C) and  $CSII_{high}$  (D) results considering all the noises (Ambulance Siren, Baby Crying, Babble and Street) for each SNR from database  $D_F$ . Each box represents the denoising performance results for each method evaluated. The color of the legend, associated to each evaluated method shown in the subfigure 16.A, refers to the same methods for all subfigures.

683 samples (Hamming window),  $F=512$  bins,  $K_N=256$  ambient noise components,  
684  $K_S=16$  biomedical components,  $M=50$  iterations of the multiplicative update  
685 rules,  $i_o=3$  incremental iterations and  $T=32$  frames per second, implies a total  
686 computational cost of  $3.32 \times 10^9$  multiplications per second and  $1.99 \times 10^9$   
687 additions per second; ii) MSS with  $N=400$  samples (Hamming window),  $F=400$   
688 bins and  $T=23$  frames per second, obtains a total computational cost of  $2.78 \times$   
689  $10^5$  multiplications per second and  $4.10 \times 10^5$  additions per second; and iii)  
690 NLMS with  $L=10$  coefficients of the adaptive filter and  $f_s=8000$  samples per  
691 second, obtains a total computational cost of  $2.48 \times 10^5$  multiplications per

692 second and  $2.40 \times 10^5$  additions per second. Results confirm that the proposed  
693 method requires a greater number of elementary operations compared to the  
694 other evaluated methods, since most of its computational cost is derived from  
695 the iterative process based on the multiplicative update rules to estimate the  
696 bases and activations matrices used to model the spectro-temporal behavior of  
697 biomedical sounds and ambient noise factorized into the NMPCF decomposition  
698 in order to extract those spectral patterns active in the magnitude spectrograms  
699 corresponding to the internal and external signal (see section 2.3.3).

Algorithms	Multiplications	Additions
Proposed method	$i_oM[FT(6 + 5K_N + 3K_S) + K_N(3T + 2F) + K_S(T + F)] + 2TN \log_2(N)$	$i_oM[FT(3K_N + 3K_S - 2) - T(2K_N + 2K_S) - FK_S] + 4TN \log_2(N)$
MSS	$13FT + 2TN \log_2(N)$	$10FT + 4TN \log_2(N)$
NLMS	$(3L + 1)f_s$	$3Lf_s$

Table 2: Assessment of computational complexity in terms of the number of operations (multiplication and addition). Note that the number of operations has been calculated from the set of parameters on which the computational cost of each evaluated algorithm depends.

#### 700 4. Conclusions and Future Work

701 In this work, we propose an incremental algorithm based on multichannel  
702 non-negative matrix partial co-factorization (NMPCF) for ambient denoising in  
703 auscultation focusing on high noisy environments with a Signal-to-Noise Ratio  
704 (SNR) lower than 0 dB. The first contribution applies NMPCF from a multi-  
705 channel point of view assuming that ambient noise can be modelled as repetitive  
706 sound events found in two single-channel audio inputs simultaneously captured  
707 by means of different recording devices. The second contribution proposes an  
708 incremental algorithm, based on the previous multichannel NMPCF, that re-  
709 fines the estimated biomedical spectrogram through a set of incremental stages  
710 by eliminating most of the ambient noises that was not removed in the previous  
711 stage.

712 The optimization process, using a database that is not adapted to any spe-  
713 cific patient's body or room, indicates that the best performance of the proposed  
714 method is obtained using a higher number of noise bases compared to the num-  
715 ber of biomedical bases. It suggests that the energy distribution of the types of  
716 ambient noises analyzed is more complex compared to biomedical sounds due  
717 to the greater spectral diversity shown by the time-frequency structures of such  
718 noises.

719 Two databases have been created in order to simulate two different acoustic  
720 scenarios depending on whether the effects of propagation on the patient's body  
721 material and the acoustics of the room are inactive (ideal scenario) or active (real  
722 scenario). Each database is composed of a broad set of biomedical sounds (heart  
723 and lung) and ambient sounds mixed together with various levels of SNR.

724 The results obtained in the ideal scenario indicate that the proposed method  
725 outperforms MSS and NLMS, in terms of SDR, in most of the levels SNR  
726 showing the most robust SDR performance mainly in high noise environments  
727 ( $\text{SNR} \in [-20 \text{ dB}, -10 \text{ dB}]$ ). Although NLMS is slightly better, in terms of SIR,  
728 compared to the proposed method, this fact occurs at the expense of losing  
729 a large amount of biomedical content. SDR and SIR improvements provided  
730 by MSS are ranked in the last position since the lowest levels SNR cause sig-  
731 nificant distortion of biomedical sounds at low spectral ranges. A remarkable  
732 advantage of the proposed method, compared to MSS and NLMS, is the high  
733 robustness of the acoustic quality of the estimated biomedical sounds when the  
734 two single-channel inputs suffer from a delay between them. As a consequence,  
735 the proposed method can be considered a more reliable ambient denoising ap-  
736 proach since results obtained assuming ambient noises as repetitive sound events  
737 found in both the internal and external spectrograms report less dependence on  
738 the ratio of ambient noise to biomedical content or any temporal misalignment  
739 between the internal and external signal as occurs with MSS and NLMS.

740 From the real scenario, results show that the proposed method significantly  
741 outperforms MSS and NLMS for each SNR evaluated. The high and stable  
742 trend of SDR and SIR displayed by the proposed method for each SNR in com-

743 parison with the rest of the methods is remarkable. Although NLMS shows  
744 better denoising performance compared to MSS, none of them can be consid-  
745 ered competitive methods since their SDR results  $\leq 1$  dB when all ambient  
746 noises are averaged so, it implies a huge loss of biomedical content despite the  
747 fact that NLMS gets a great isolated biomedical signal. Comparing the denois-  
748 ing performance between the ideal and real scenario, it can be observed a high  
749 robustness of the proposed method compared to MSS and NLMS. Specifically,  
750 the average reduction of the SDR and SIR improvement suffered by the pro-  
751 posed averaging all ambient noises and taking into account the effects of the  
752 propagation of the patient’s body material and the acoustics of the room is  $\leq$   
753 1 dB and  $\leq 5$  dB while the same SDR and SIR reduction increases to 14 dB  
754 and 15 dB for MSS, and 18 dB and 21 dB for NLMS. Moreover, another set  
755 of metrics (NCM, CSII<sub>low</sub>, CSII<sub>med</sub> and CSII<sub>high</sub>) confirm that the proposed  
756 method provides the best ambient denoising performance by means of maximiz-  
757 ing the ambient denoising at the expense of minimizing the loss of biomedical  
758 energy. Again, these results demonstrate that the multichannel co-factorization  
759 approach achieves a more adequate behavior in the removal of ambient noise  
760 compared to the previous spectral subtraction or adaptive filtering approaches.  
761 As a drawback, highlight the high computational cost of the proposed method  
762 mainly to the number of operation applied to the multiplicative update rules of  
763 the co-factorization process.

764 Future work will focus on two directions: (i) novel algorithms applied to the  
765 removal of some of the most disturbing acoustically active ambient noises in clin-  
766 ical emergency situations, such as noise inside a helicopter or ambulance when  
767 urgent monitoring is required, and, (ii) real-time multichannel non-negative ma-  
768 trix partial co-factorization approaches for biomedical ambient denoising.

## 769 Funding

770 This work was supported by the Programa Operativo FEDER Andalucia  
771 2014-2020 under project with reference 1257914 and the Ministry of Economy,



772 Knowledge and University, Junta de Andalucia under Project P18-RT-1994.

### 773 **Acknowledgment**

774 The authors would like to thank Dr. Dinko Oletic and Dr. Vedran Bilas  
775 for sharing their lung recordings. The authors would like to thank the pulmo-  
776 nologist Gerardo Perez Chica from the University Hospital of Jaen (Spain) for  
777 his assistance related to ambient noises. Finally, we would like to thank the  
778 anonymous reviewers for their helpful and constructive comments that greatly  
779 contributed to improve the final version of the paper.

### 780 **References**

- 781 [1] A. K. Abbas, R. Bassam, Phonocardiography signal processing, Synthesis  
782 Lectures on Biomedical Engineering 4 (1) (2009) 1–194.
- 783 [2] M. Sarkar, I. Madabhavi, N. Niranjana, M. Dogra, Auscultation of the res-  
784 piratory system, *Annals of Thoracic Medicine* 10 (3) (2015) 158–168.
- 785 [3] S. Taplidou, L. Hadjileontiadis, Wheeze detection based on time-frequency  
786 analysis of breath sounds, *Computers in biology and medicine* 37 (8) (2007)  
787 1073–1083.
- 788 [4] D. Kumar, P. Carvalho, M. Antunes, J. Henriques, Noise detection during  
789 heart sound recording, in: *31st Annual International Conference of the*  
790 *IEEE Engineering in Medicine and Biology Society, IEEE, 2009*, pp. 3119–  
791 3123.
- 792 [5] T. Tsalaile, S. Sanei, Separation of heart sound signal from lung sound  
793 signal by adaptive line enhancement, in: *15th European Signal Processing*  
794 *Conference, IEEE, 2007*, pp. 1231–1235.
- 795 [6] C. Lin, E. Hasting, Blind source separation of heart and lung sounds based  
796 on nonnegative matrix factorization, in: *International symposium on in-*  
797 *telligent signal processing and communication systems (ISPACS), IEEE,*  
798 *2013*, pp. 731–736.

- 799 [7] F. Canadas-Quesada, N. Ruiz-Reyes, J. Carabias-Orti, P. Vera-Candeas,  
800 J. Fuertes-Garcia, A non-negative matrix factorization approach based on  
801 spectro-temporal clustering to extract heart sounds, *Applied Acoustics* 125  
802 (2017) 7–19.
- 803 [8] G. Serbes, C. Sakar, Y. Kahya, N. Aydin, Pulmonary crackle detection  
804 using time–frequency and time–scale analysis, *Digital Signal Processing*  
805 23 (3) (2013) 1012–1021.
- 806 [9] M. Zivanovic, M. Gonzalez-Izal, Quasi-periodic modeling for heart sound  
807 localization and suppression in lung sounds, *Biomedical Signal Processing*  
808 *Control* 8 (2013) 586–595.
- 809 [10] V. Varghees, K. Ramachandran, A novel heart sound activity detection  
810 framework for automated heart sound analysis, *Biomed. Signal Process.*  
811 *Control.* 13 (2014) 174–188.
- 812 [11] J. Torre-Cruz, F. Canadas-Quesada, J. Carabias-Orti, P. Vera-Candeas,  
813 N. Ruiz-Reyes, A novel wheezing detection approach based on constrained  
814 non-negative matrix factorization, *Applied Acoustics* 148 (2019) 276–288.
- 815 [12] F. Jin, F. Sattar, D. Goh, New approaches for spectro-temporal feature  
816 extraction with applications to respiratory sound classification, *Neurocom-*  
817 *puting* 123 (2014) 362–271.
- 818 [13] S. Raj, K. Ray, O. Shankar, Cardiac arrhythmia beat classification using  
819 dost and pso tuned svm, *Computer methods and programs in biomedicine*  
820 136 (2016) 163–77.
- 821 [14] P. Li, Y. Wang, J. He, L. Wang, Y. Tian, T.-s. Zhou, T. Li, J.-s. Li, High-  
822 performance personalized heartbeat classification model for long-term ecg  
823 signal, *IEEE Transactions on Biomedical Engineering* 64 (1) (2016) 78–86.
- 824 [15] D. Bardou, K. Zhang, S. Ahmad, Lung sounds classification using convolu-  
825 tional neural networks, *Artificial Intelligence in Medicine* 88 (2018) 58–69.

- 826 [16] R. X. A. Pramono, S. A. Imtiaz, E. Rodriguez-Villegas, Evaluation of fea-  
827 tures for classification of wheezes and normal respiratory sounds., *PloS one*  
828 14 (3) (2019) e0213659–e0213659.
- 829 [17] A. Suzuki, C. Sumi, K. Nakayama, M. Mori, Real-time adaptive cancelling  
830 of ambient noise in lung sound measurement, *Medical and Biological En-*  
831 *gineering and Computing* 33 (5) (1995) 704–708.
- 832 [18] S. B. Patel, T. F. Callahan, M. G. Callahan, J. T. Jones, G. P. Graber, K. S.  
833 Foster, K. Glifort, G. R. Wodicka, An adaptive noise reduction stethoscope  
834 for auscultation in high noise environments, *The Journal of the Acoustical*  
835 *Society of America* 103 (5) (1998) 2483–2491.
- 836 [19] J. S. Fleeter, G. R. Wodicka, Auscultation of heart and lung sounds in high-  
837 noise environments using adaptive filters, *The Journal of the Acoustical*  
838 *Society of America* 104 (3) (1998) 1781–1781.
- 839 [20] D. Della Giustina, M. Riva, F. Belloni, M. Malcangi, Embedding a mul-  
840 tichannel environmental noise cancellation algorithm into an electronic  
841 stethoscope, *International Journal of Circuits/System and Signal Process-*  
842 *ing* (2) (2011).
- 843 [21] G. Nelson, R. Rajamani, A. Erdman, Noise control challenges for auscultation  
844 on medical evacuation helicopters, *Applied Acoustics* 80 (2014) 68–78.
- 845 [22] N. W. Evans, J. S. Mason, W.-M. Liu, B. Fauve, An assessment on the fun-  
846 damental limitations of spectral subtraction, in: *2006 IEEE International*  
847 *Conference on Acoustics Speech and Signal Processing Proceedings*, Vol. 1,  
848 IEEE, 2006, pp. I–I.
- 849 [23] G. Chang, Y. Lai, Performance evaluation and enhancement of lung sound  
850 recognition system in two real noisy environments, *Computer methods and*  
851 *programs in biomedicine* 97 (2) (2015) 141–150.
- 852 [24] D. Emmanouilidou, E. McCollum, D. Park, M. Elhilali, Adaptive noise  
853 suppression of pediatric lung auscultations with real applications to noisy

- 854 clinical settings in developing countries, *IEEE Transactions on Biomedical*  
855 *Engineering* 62 (9) (2015) 2279–2288.
- 856 [25] D. Emmanouilidou, E. D. McCollum, D. E. Park, M. Elhilali, Computerized  
857 lung sound screening for pediatric auscultation in noisy field environments,  
858 *IEEE Transactions on Biomedical Engineering* 65 (7) (2018) 1564–1574.
- 859 [26] Y. Hu, G. Liu, Separation of singing voice using nonnegative matrix par-  
860 tial co-factorization for singer identification, *IEEE/ACM Transactions on*  
861 *Audio, Speech, and Language Processing* 23 (4) (2015) 643–653.
- 862 [27] J. Yoo, M. Kim, K. Kang, S. Choi, Nonnegative matrix partial co-  
863 factorization for drum source separation, in: *2010 IEEE International Con-*  
864 *ference on Acoustics, Speech and Signal Processing*, IEEE, 2010, pp. 1942–  
865 1945.
- 866 [28] M. Kim, J. Yoo, K. Kang, S. Choi, Blind rhythmic source separation:  
867 Nonnegativity and repeatability, in: *2010 IEEE International Conference*  
868 *on Acoustics, Speech and Signal Processing*, IEEE, 2010, pp. 2006–2009.
- 869 [29] M. Kim, J. Yoo, K. Kang, S. Choi, Nonnegative matrix partial co-  
870 factorization for spec- tral and temporal drum source separation, *IEEE*  
871 *Journal on Selected Topics in Signal Processing* 5 (6) (2011) 1192–1204.
- 872 [30] N. Seichepine, S. Essid, C. Févotte, O. Cappé, Soft nonnegative matrix  
873 co-factorization, *IEEE Transactions on Signal Processing* 62 (22) (2014)  
874 5940–5949.
- 875 [31] J. De La Torre Cruz, F. J. Cañadas Quesada, N. Ruiz Reyes, P. Vera Can-  
876 deas, J. J. Carabias Orti, Wheezing sound separation based on informed  
877 inter-segment non-negative matrix partial co-factorization, *Sensors* 20 (9)  
878 (2020) 2679.
- 879 [32] D. Badawy, N. Duong, A. Ozerov, On-the-fly audio source separation—a  
880 novel user-friendly framework, *IEEE/ACM Transactions on Audio, Speech,*  
881 *and Language* 25 (2) (2016) 261–272.

- 882 [33] V. Bisot, R. Serizel, S. Essid, G. Richard, Leveraging deep neural networks  
883 with nonnegative representations for improved environmental sound classi-  
884 fication, in: IEEE International Workshop on Machine Learning for Signal  
885 Processing (MLSP), IEEE, 2017, pp. 1–6.
- 886 [34] A. Mesaros, A. Diment, B. Elizalde, T. Heittola, E. Vincent, R. Bhik-  
887 sha, T. Virtanen, Sound event detection in the dcase 2017 challenge,  
888 IEEE/ACM Transactions on Audio, Speech and Language Processing  
889 27 (6) (2019) 992–1006.
- 890 [35] Freesound by Music Technology Group, Universitat Pompeu Fabra, <https://freesound.org/>,  
891 online. Accessed: 2020-04-27 (2005).
- 892 [36] Findsound by Comparisonics Corporation, <https://www.findsounds.com/>,  
893 online. Accessed: 2020-04-27 (2020).
- 894 [37] Detection and Classification of Acoustic Scenes and Events DCASE  
895 2017 Challenge. Detection of rare sound events (Tampere Uni-  
896 versity of Technology), [http://www.cs.tut.fi/sgn/arg/dcase2017/  
897 challenge/task-rare-sound-event-detection](http://www.cs.tut.fi/sgn/arg/dcase2017/challenge/task-rare-sound-event-detection), online. Accessed: 2020-  
898 04-27 (2017).
- 899 [38] Signal Processing Information Base (SPIB). NOISEX database. Speech  
900 Babble, <http://spib.linse.ufsc.br/noise.html>,  
901 online. Accessed: 2020-04-27 (1990).
- 902 [39] ETSI TS 103 224 V1. Speech and multimedia Transmission Quality  
903 (STQ); A sound field reproduction method for terminal testing includ-  
904 ing a background noise database. Background Noise Database: cafe-  
905 teria and pub, [https://docbox.etsi.org/stq/Open/TS%20103%20224%  
906 20Background%20Noise%20Database/Binaural](https://docbox.etsi.org/stq/Open/TS%20103%20224%20Background%20Noise%20Database/Binaural), online. Accessed: 2020-  
907 04-27 (2014).
- 908 [40] Detection and Classification of Acoustic Scenes and Events DCASE  
909 2017 Challenge. Sound event detection in real life audio (Tampere

- 910 University of Technology), [http://www.cs.tut.fi/sgn/arg/dcase2017/  
911 challenge/task-acoustic-scene-classification](http://www.cs.tut.fi/sgn/arg/dcase2017/challenge/task-acoustic-scene-classification), online. Accessed:  
912 2020-04-27 (2017).
- 913 [41] TUT Sound events 2017, Development dataset, [https://zenodo.org/  
914 record/814831](https://zenodo.org/record/814831), online. Accessed: 2020-04-27 (2017).
- 915 [42] TUT Sound events 2017, Evaluation dataset, [https://zenodo.org/  
916 record/1040179](https://zenodo.org/record/1040179), online. Accessed: 2020-04-27 (2017).
- 917 [43] PASCAL Classifying heart sounds challenge, [http://www.  
918 peterjbentley.com/heartchallenge/](http://www.peterjbentley.com/heartchallenge/), online. Accessed: 2020-04-27  
919 (2011).
- 920 [44] PhysioNet/CinC challenge. National Institute of General Medical Sciences  
921 and the National Institute of Biomedical Imaging and Bioengineering,  
922 <https://www.physionet.org/physiobank/database/challenge/2016/>,  
923 online. Accessed: 2020-04-27 (2013).
- 924 [45] S. Charleston-Villalobos, L. Dominguez-Robert, R. Gonzalez-Camarena,  
925 A. Aljama-Corrales, Heart sounds interference cancellation in lung sounds,  
926 in: 2006 International Conference of the IEEE Engineering in Medicine and  
927 Biology Society, IEEE, 2006, pp. 1694–1697.
- 928 [46] S. M. Debbal, F. Bereksi-Reguig, Spectral analysis of the pcg signals, The  
929 Internet journal of microbiology 2 (2006).
- 930 [47] D. Oletic, V. Bilas, Asthmatic wheeze detection from compressively sensed  
931 respiratory sound spectra, IEEE journal of biomedical and health infor-  
932 matics 22 (5) (2018) 1406–1414.
- 933 [48] A. Sovijarvi, J. Vanderschoot, J. Earis, Standardization of computerized  
934 respiratory sound analysis, European Respiratory Review 10 (77) (2000)  
935 585–585.

- 936 [49] S. Reichert, R. Gass, C. Brandt, E. Andrès, Analysis of respiratory sounds:  
937 state of the art, *Clinical medicine. Circulatory, respiratory and pulmonary*  
938 *medicine* 2 (2008) CCRPM-S530.
- 939 [50] S. Haykin, *Adaptive filter theory*, Englewood Cliffs, NJ: PrenticeHall, 1996.
- 940 [51] DSP System Toolbox, Filter Implementation, Adaptive Filters, [https://es.mathworks.com/help/dsp/ref/dsp.lmsfilter-system-object.](https://es.mathworks.com/help/dsp/ref/dsp.lmsfilter-system-object.html)  
941 [html](https://es.mathworks.com/help/dsp/ref/dsp.lmsfilter-system-object.html).  
942
- 943 [52] J. Torre-Cruz, F. Canadas-Quesada, S. García-Galán, N. Ruiz-Reyes,  
944 P. Vera-Candeas, J. Carabias-Orti, A constrained tonal semi-supervised  
945 non-negative matrix factorization to classify presence/absence of wheezing  
946 in respiratory sounds, *Applied Acoustics* 161 (2020) 107–188.
- 947 [53] E. Vincent, R. Gribonval, C. Févotte, Performance measurement in blind  
948 audio source separation, *IEEE transactions on audio, speech, and language*  
949 *processing* 14 (4) (2006) 1462–1469.
- 950 [54] C. Févotte, R. Gribonval, E. Vincent, *Bss\_eval toolbox user guide–revision*  
951 *2.0* (2005).
- 952 [55] Y. Matsui, S. Makino, N. Ono, T. Yamada, Multiple far noise suppression  
953 in a real environment using transfer-function-gain nmf, in: *2017 25th Eu-*  
954 *ropean Signal Processing Conference (EUSIPCO)*, IEEE, 2017, pp. 2314–  
955 2318.
- 956 [56] A. Liutkus, D. Fitzgerald, Z. Rafii, Scalable audio separation with light  
957 kernel additive modelling, in: *IEEE International Conference on Acoustics,*  
958 *Speech and Signal Processing (ICASSP)*, IEEE, 2015, pp. 76–80.
- 959 [57] F. Li, M. Akagi, Blind monaural singing voice separation using rank-1  
960 constraint robust principal component analysis and vocal activity detection,  
961 *Neurocomputing* 350 (2019) 44–52.

- 962 [58] S. Venkataramani, C. Subakan, P. Smaragdis, Neural network alternatives  
963 toconvolutive audio models for source separation, in: IEEE International  
964 Workshop on Machine Learning for Signal Processing, IEEE, 2017, pp. 1–6.
- 965 [59] P. C. Loizou, Speech enhancement: theory and practice, CRC press, 2013.
- 966 [60] G.-C. Chang, A comparative analysis of various respiratory sound denois-  
967 ing methods, in: 2016 International Conference on Machine Learning and  
968 Cybernetics (ICMLC), Vol. 2, IEEE, 2016, pp. 514–518.
- 969 [61] S. Leng, R. San Tan, K. Tshun, C. Chai, C. Wang, G. D., L. Zhong,  
970 The electronic stethoscope, Biomedical engineering online 66 (2015). doi:  
971 10.1186/s12938-015-0056-y.
- 972 [62] R. M. Corey, N. Tsuda, A. C. Singer, Wearable microphone impulse re-  
973 sponses (2018). doi:10.13012/B2IDB-1932389\_v1.
- 974 [63] R. M. Corey, N. Tsuda, A. C. Singer, Acoustic impulse responses for wear-  
975 able audio devices, in: ICASSP 2019 - 2019 IEEE International Conference  
976 on Acoustics, Speech and Signal Processing (ICASSP), 2019, pp. 216–220.  
977 doi:10.1109/ICASSP.2019.8682733.
- 978 [64] J. B. Allen, D. A. Berkley, Image method for efficiently simulating small-  
979 room acoustics, The Journal of the Acoustical Society of America 65 (4)  
980 (1979) 943–950.