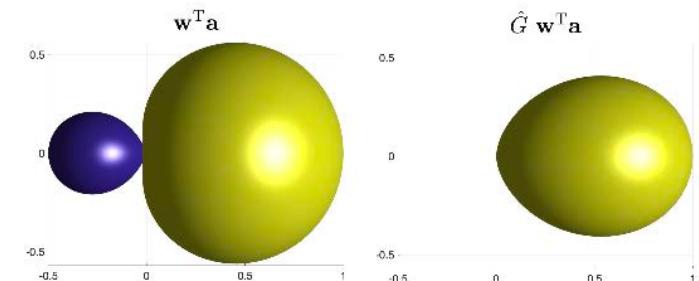
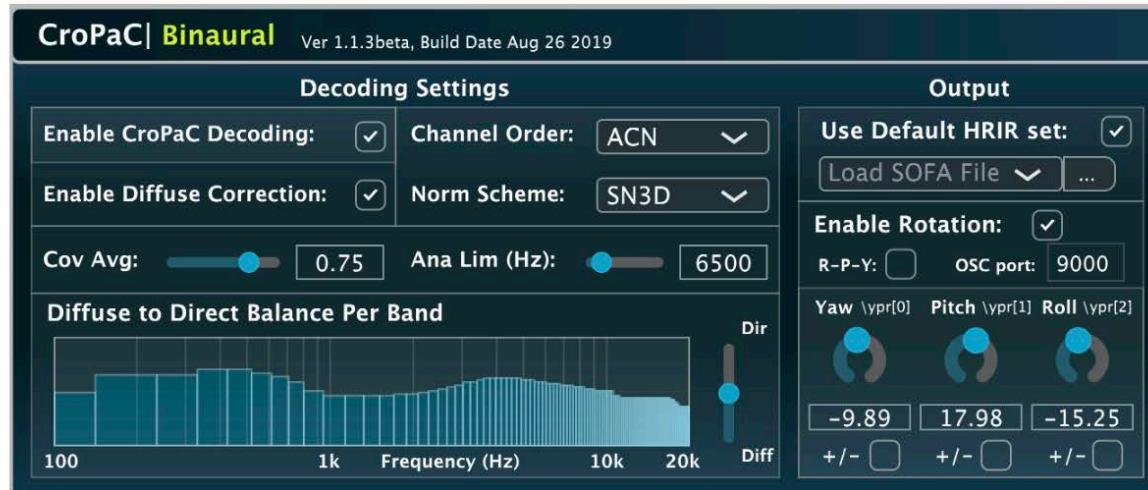




Bernard Lagnel  
Novembre 2022

## Description du plug-in



Un formateur de faisceau hyper cardioïde de premier ordre sans (à gauche) et avec (à droite) le post-filtre CroPaC...

### Décodeur ambisonique/binaural paramétrique...

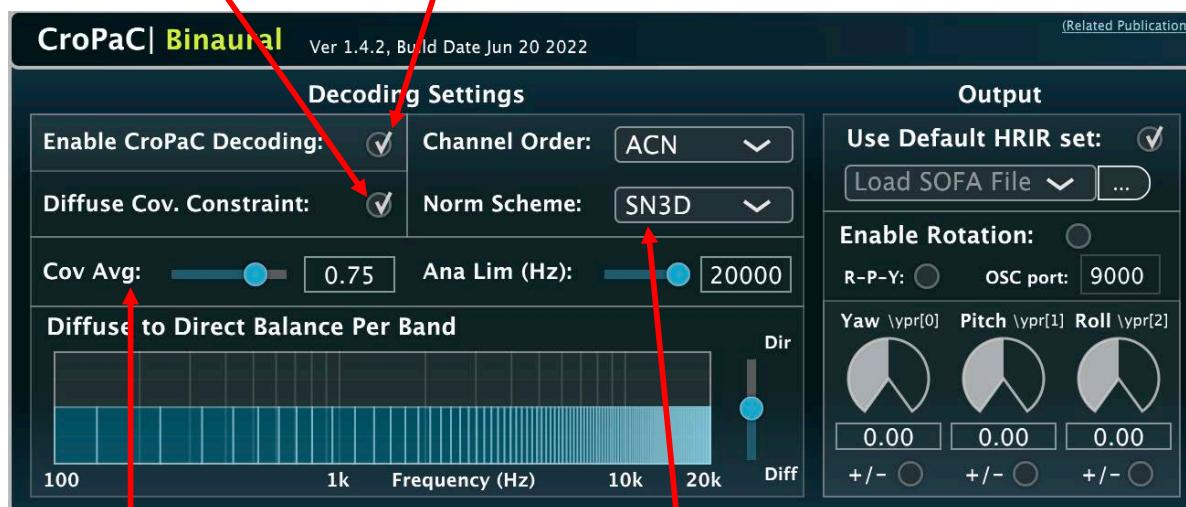
Il est spécifiquement destiné à reproduire une entrée de premier ordre avec une précision perceptuelle élevée. La méthode proposée suppose un modèle de champ sonore d'une source et d'une composante ambiante non isotrope par bande étroite. Il utilise ensuite le post-filtre Cross-Pattern Cohérence (CroPaC), afin de séparer ces composants avec une sélectivité spatiale améliorée. Les résultats des tests d'écoute indiquent que la méthode proposée lorsqu'elle utilise une entrée de premier ordre, fonctionne de manière similaire à la reproduction Ambisonics de troisième ordre !!

Le plug-in permet à l'utilisateur d'importer ses propres HRTF via des fichiers SOFA et prend également en charge le suivi de la tête via des messages OSC. L'utilisateur peut également influer sur l'équilibre direct/diffus par bande de fréquence ; notez que les flux sont équilibrés lorsqu'ils sont définis au milieu (par défaut).

<https://leomccormack.github.io/sparta-site/docs/plugins/cropac-binaural/>

Active/désactive la contrainte de covariance diffuse appliquée à la matrice de décodage. C'est la partie 'C' du décodeur 'TAC'. Notez que ce n'est pas la même chose que d'appliquer un égaliseur à champ diffus sur le les HRIR ; il s'agit principalement d'une manipulation "spatiale", pas d'une manipulation timbrale. Notez également que, même si cela peut rendre les enregistrements sonores plus large / plus large aux ordres inférieurs, il le fait au prix de endommageant considérablement les propriétés spatiales de l'enregistrement (tout tirer sur les côtés : presque stéréo-élargissement) ; par conséquent, nous dirions qu'il n'est pas "correct" d'activer cette par défaut ... bien que cela puisse sembler assez bon dans certains cas.

Active / désactive le rendu paramétrique. Lorsqu'il est désactivé, le plugin produit un son décodé ambisonique à l'aide du décodeur MagLS.

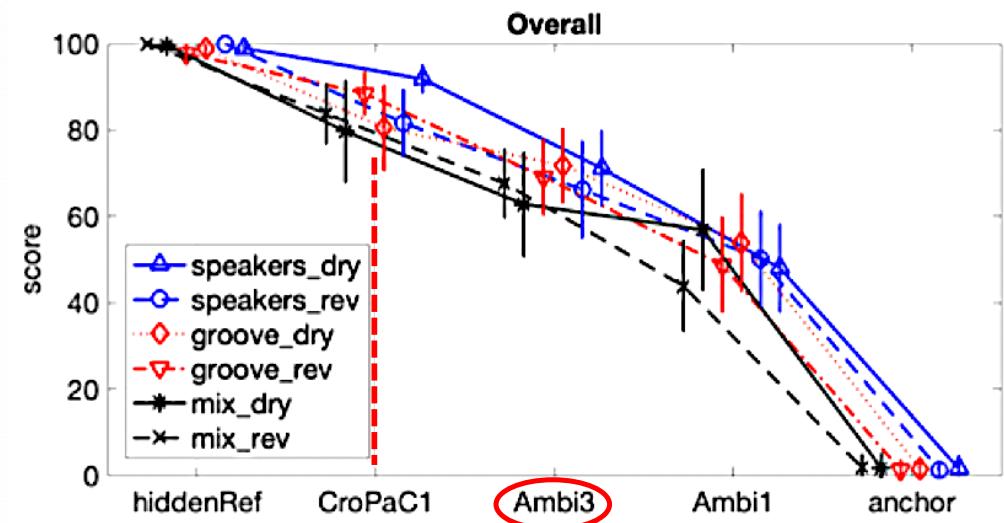
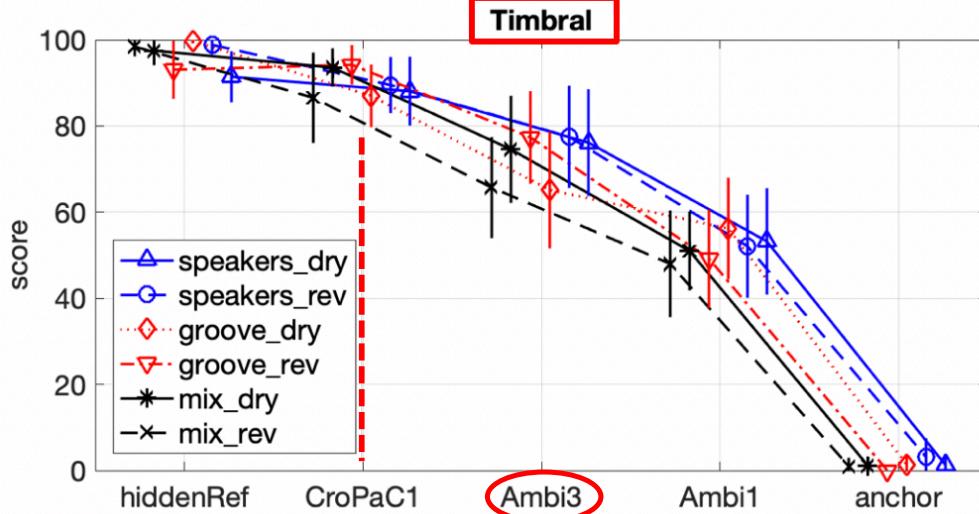
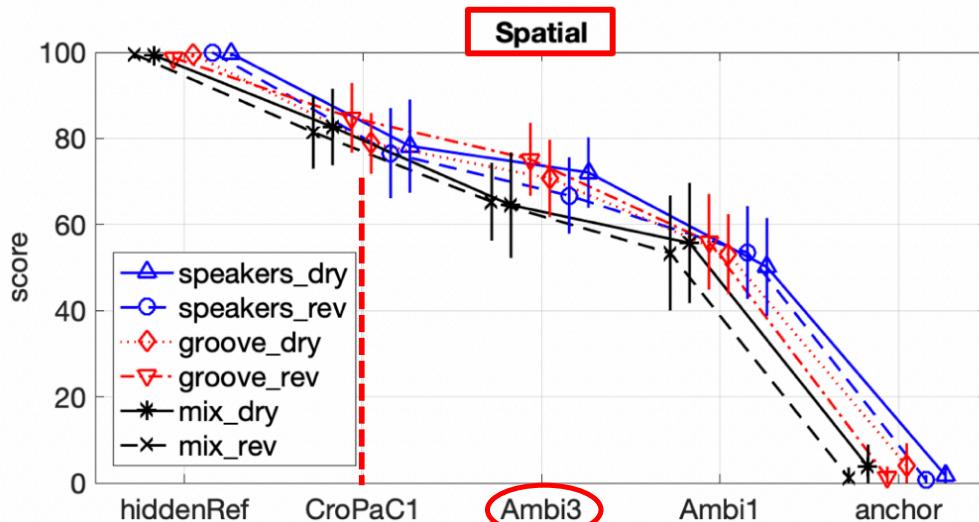


Coefficient de moyenne de la matrice de covariance (unipolaire).

Schéma de normalisation ambisonique (Notez qu' AmbiX : ACN/SN3D) .

# Résultats des tests d'écoute en un coup d'œil

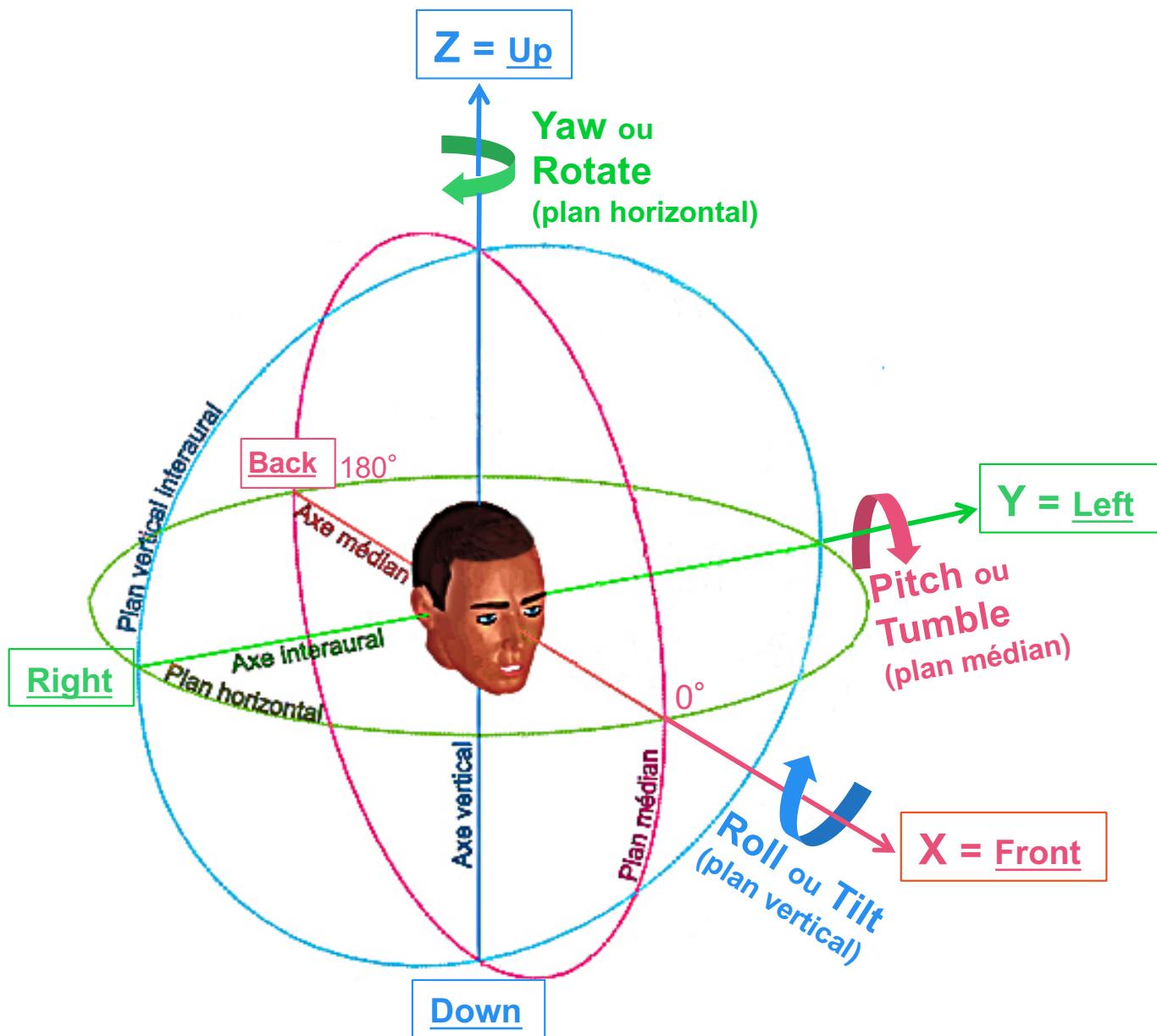
Des tests d'écoute formels indiquent que le décodeur de premier ordre proposé (CroPaC1) fonctionne de manière similaire (ou dépasse) le décodage ambisonique de rééchantillonnage spatial de troisième ordre (Ambi3), en termes d'attributs spatiaux et timbraux perçus de la reproduction ; comme indiqué dans les graphiques ci-dessous :



Il convient de souligner que l'ambisonique de troisième ordre utilise quatre fois plus de canaux d'entrée que celui de la méthode proposée !!

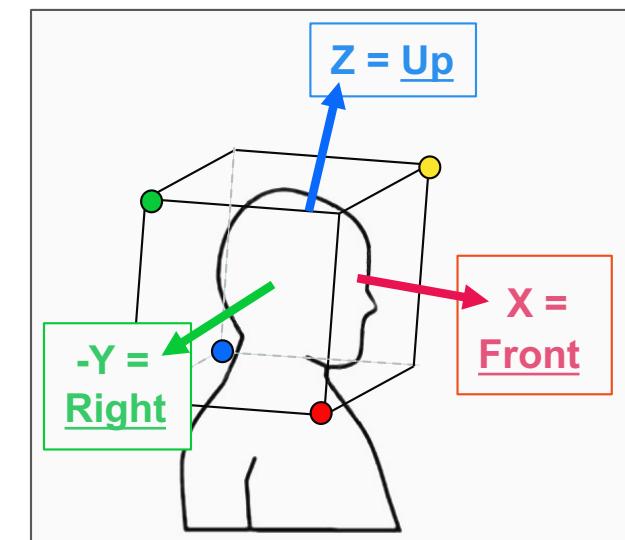
<https://leomccormack.github.io/sparta-site/docs/plugins/cropac-binaural/>

# Rotation Ambisonic 3D



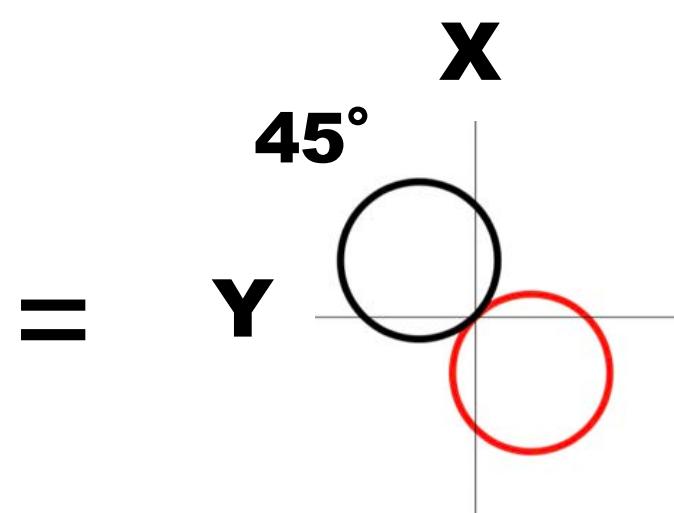
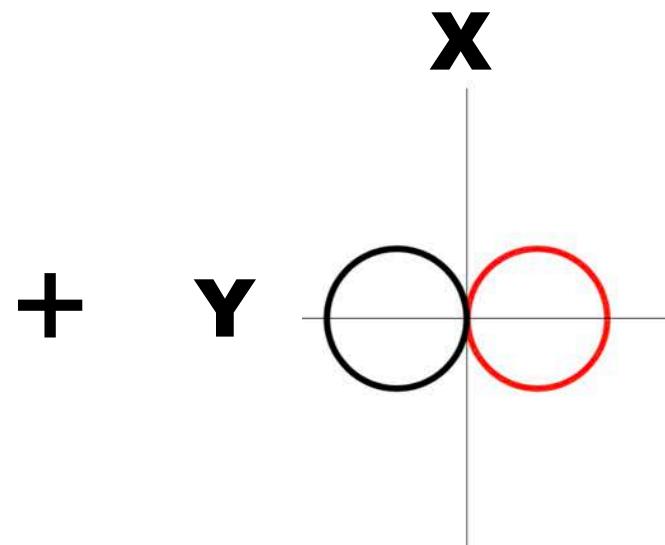
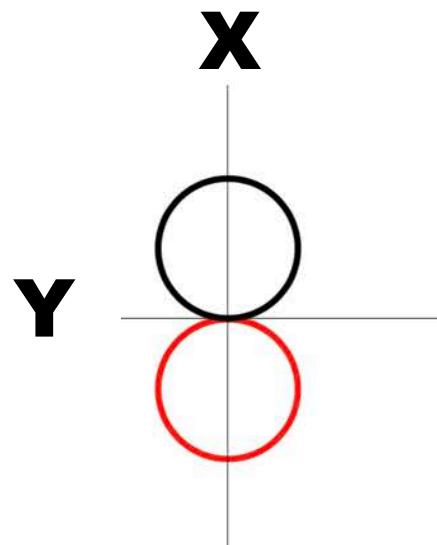
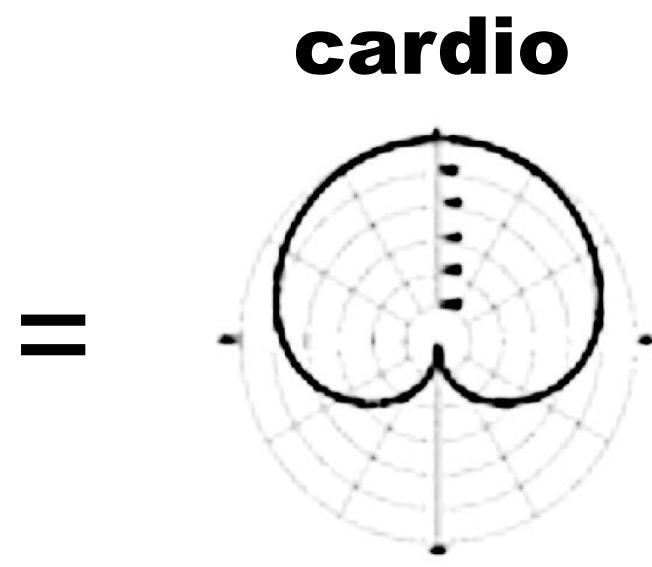
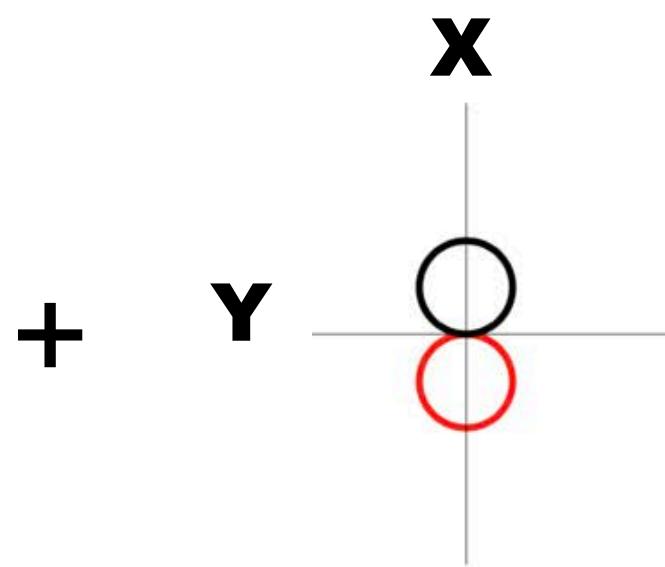
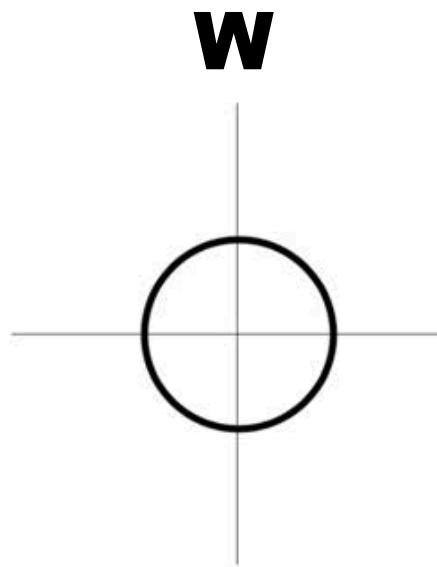
## Les 3 Plans :

1. **Plan médian :** Pitch ou Tumble
2. **Plan horizontal ou azimuthal :** Yaw ou Rotate
3. **Plan vertical ou interaural :** Roll ou Tilt



Représentation des capsules par rapport aux axes XYZ...

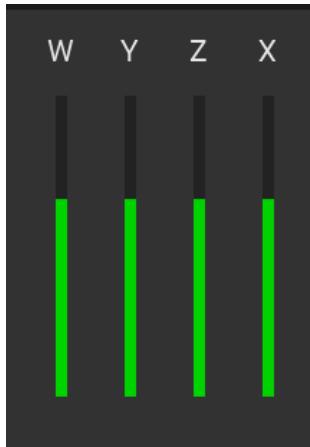
# ***Directivité et Direction***



# Sons Techniques Ambisoniques

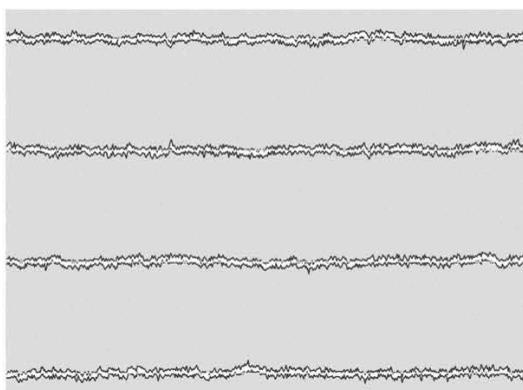
<https://www.lesonbinaural.fr/>

## Format « ambiX »



« **FuMa** » signifie «Furse-Malham», c'est à dire que l'ordre des canaux est (W, X, Y, Z) avec le canal W normalisé :  $1/\sqrt{2} = -3$  dB.

« **ambiX** » signifie l'ordre des canaux ACN avec la normalisation SN3D, c'est à dire que l'ordre des canaux est (W, Y, Z, X) sans mise à l'échelle des canaux (au même niveau).



### Bruit Rose sur 4 Pistes ©

Bruit Rose sur 4 pistes destiné au Multicanal en Quad et à l'Ambisonique (courbe de réponse, équilibre, filtre...)

Dé-corrélation + 0,0 : de 0 s à 40 s  
Corrélation + 0,25 : de 1 mn à 1 mn 40 s  
Corrélation + 0,5 : de 2 mn à 2 mn 40 s  
Corrélation + 0,75 : de 3 mn à 3 mn 40 s  
Corrélation + 1,0 : de 4 mn à 4 mn 40 s

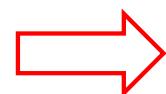
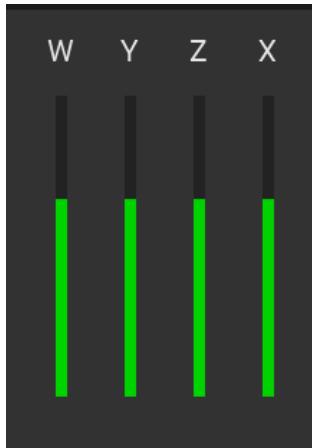
Attention au niveau -12 dBFS, coupe bas à 30 Hz.

prendre à 4 mn  
[Télécharger](#)

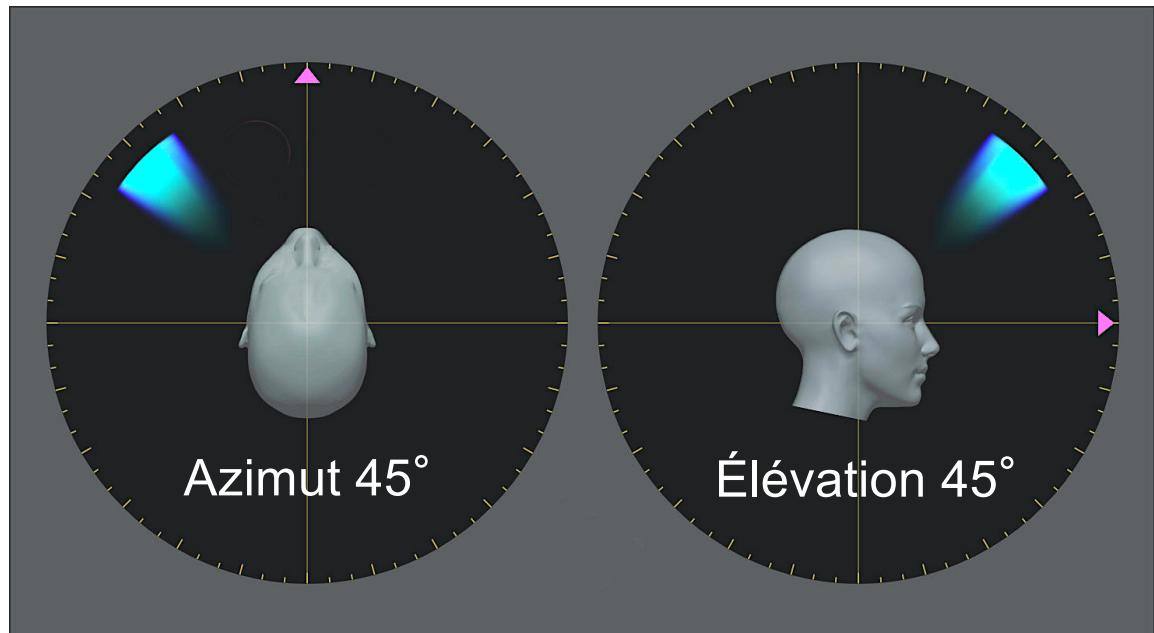
4 min 40 sec  
Quad 4.0  
L R Ls Rs  
En .WAV  
24 Bit / 48 KHz

# Direction

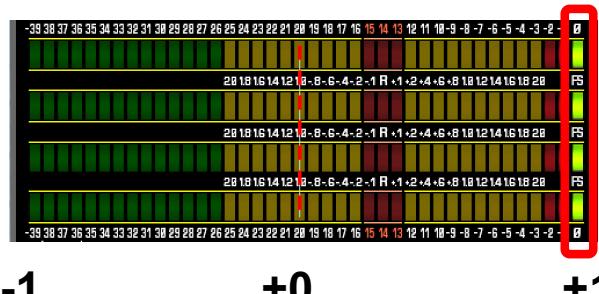
**IN : Format-B « AmbiX »**



**OUT : Binaural**



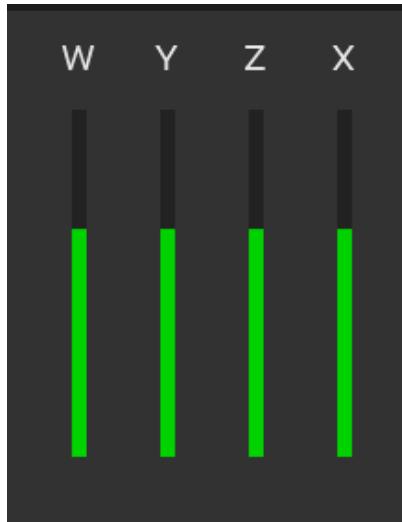
Bruit Rose Corrélé sur les 4 canaux  
(Phase à + 1 = mêmes signaux)



Télécharger les HRTF en .sofa :  
<https://www.lesonbinaural.fr/EDIT/HRTF>

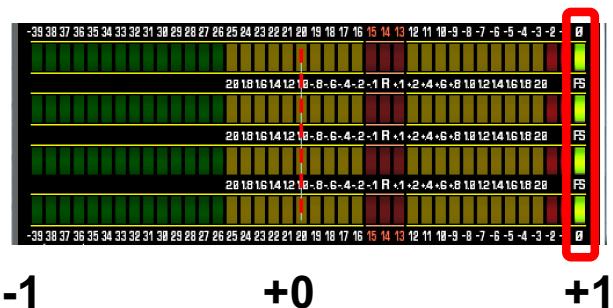
Name	Last modified	Size	Description
<a href="#">Parent Directory</a>		-	
<a href="#">BBC_RD_32POS_48K.sofa</a>	2022-11-04 18:55	4.9M	
<a href="#">SADIE_D1_48K_24bit_2...&gt;</a>	2022-11-04 18:55	35M	
<a href="#">SADIE_KU100_DFC_256_...&gt;</a>	2022-11-04 18:55	12M	
<a href="#">SENNHEISER_ORBIT.sofa</a>	2022-11-04 18:55	2.6M	
<a href="#">TH_KOLN_HRIR_FULL2DE...&gt;</a>	2022-11-04 18:55	19M	

Pour : 45° Azimuth  
45° Elévation

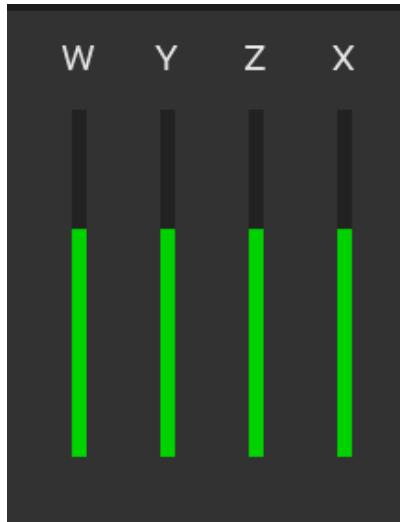


**IN : Format B AmbiX**

Bruit Rose Corrélé sur les 4 canaux  
(Phase à + 1 = mêmes signaux)

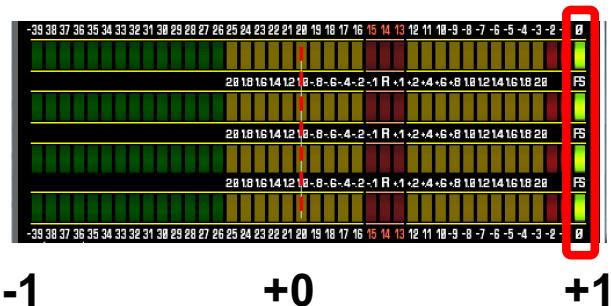


Pour : 45° Azimuth  
45° Elévation



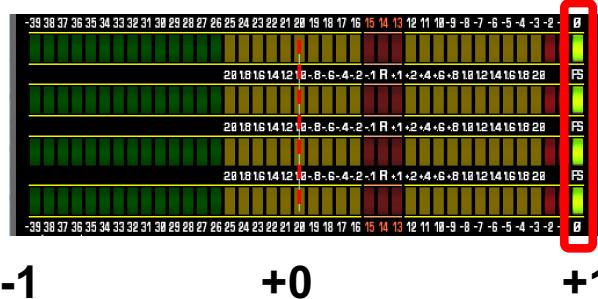
**IN : Format B AmbiX**

Bruit Rose Corrélu sur les 4 canaux  
(Phase à + 1 = mêmes signaux)



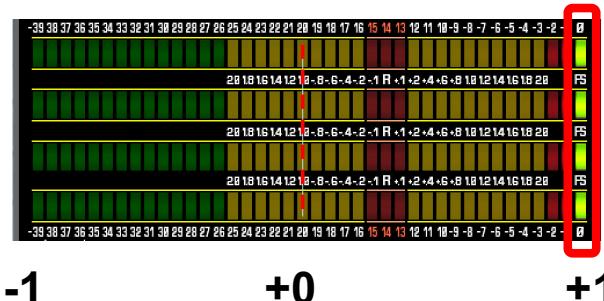
## ***IN : Format B AmbiX***

**Bruit Rose Corrélu sur les 4 canaux  
(Phase à + 1 = mêmes signaux)**



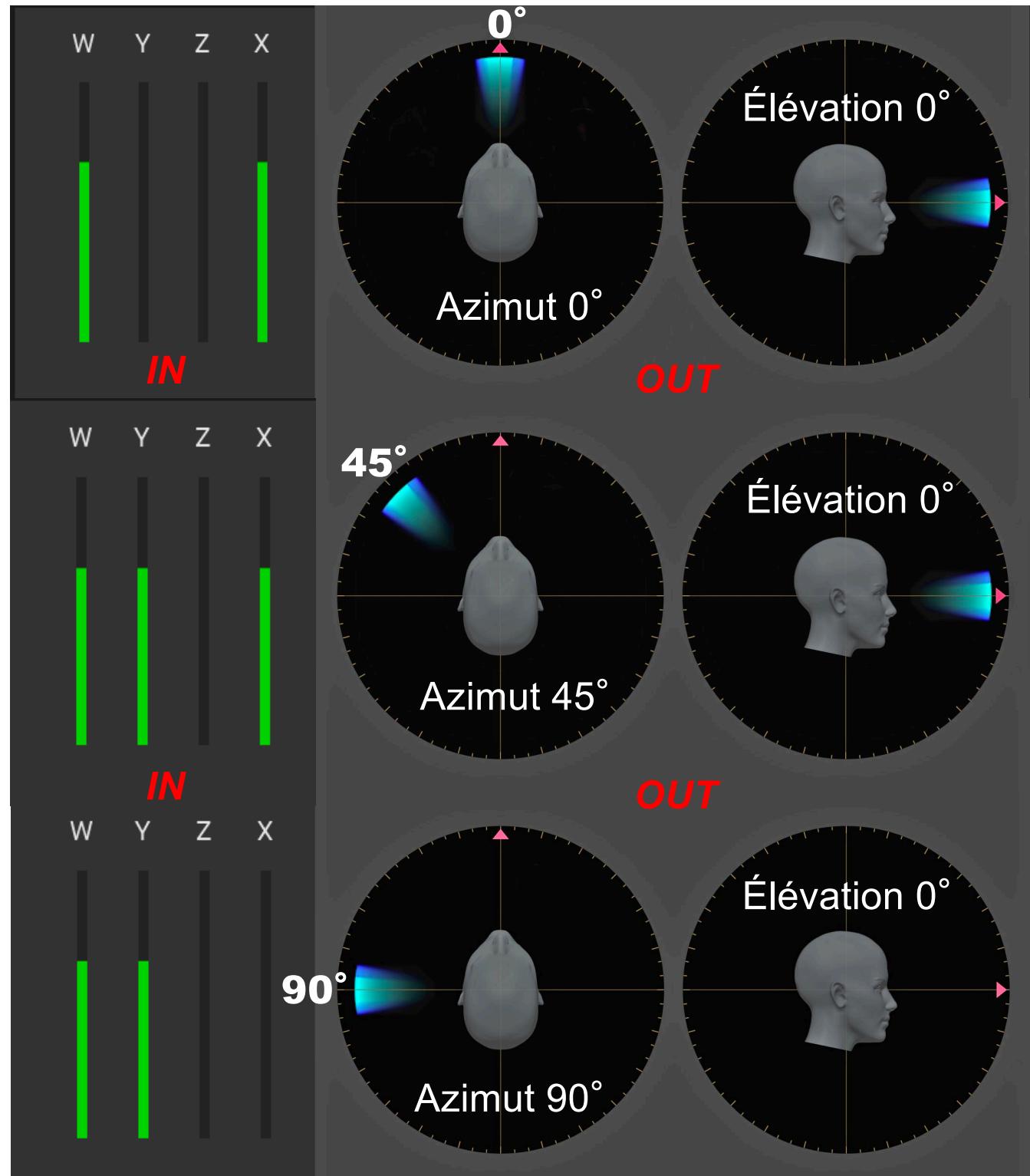
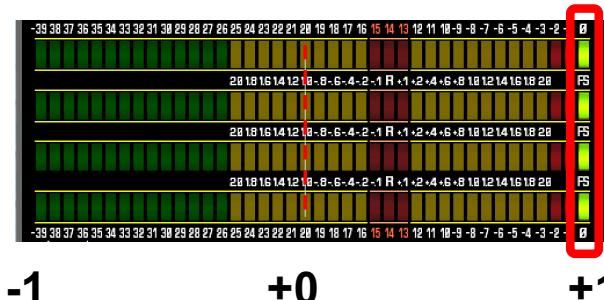
# *IN : Format B AmbiX*

**Bruit Rose Corrélu sur les 4 canaux  
(Phase à + 1 = mêmes signaux)**



## ***IN : Format B AmbiX***

**Bruit Rose Corrélé sur les 4 canaux  
(Phase à + 1 = mêmes signaux)**

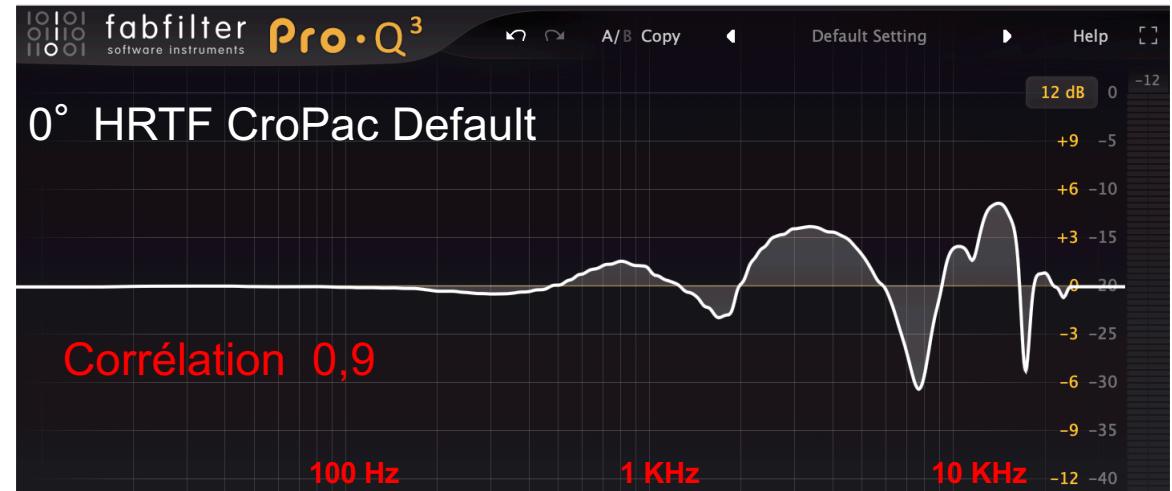
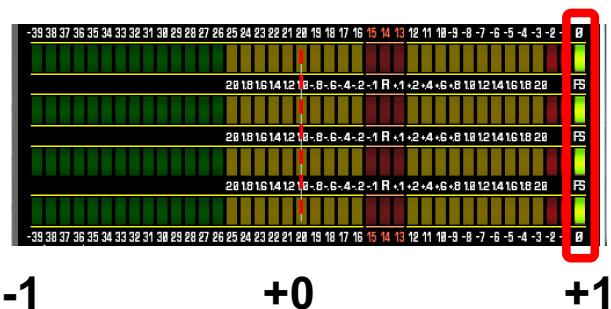


Pour : 0° Azimuth  
0° Elévation

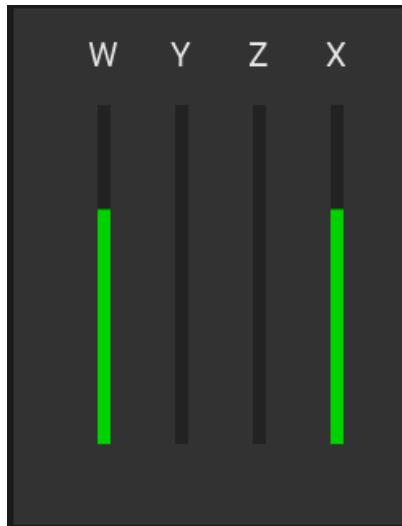


**IN : Format B AmbiX**

Bruit Rose Corrélué sur les 4 canaux  
(Phase à + 1 = mêmes signaux)

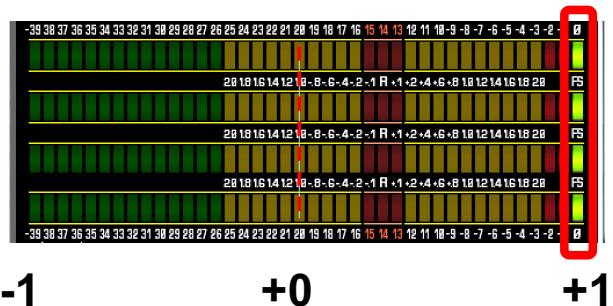


Pour : 0° Azimuth  
0° Elévation

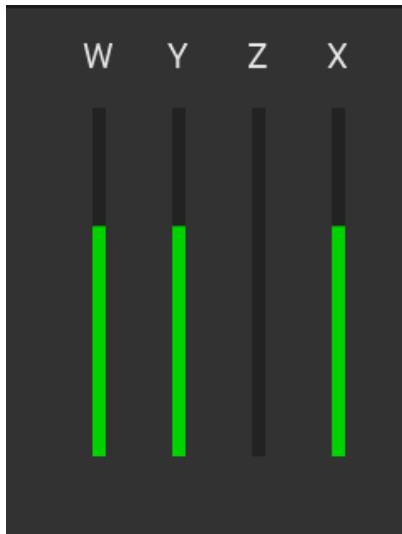


**IN : Format B AmbiX**

Bruit Rose Corrélué sur les 4 canaux  
(Phase à + 1 = mêmes signaux)

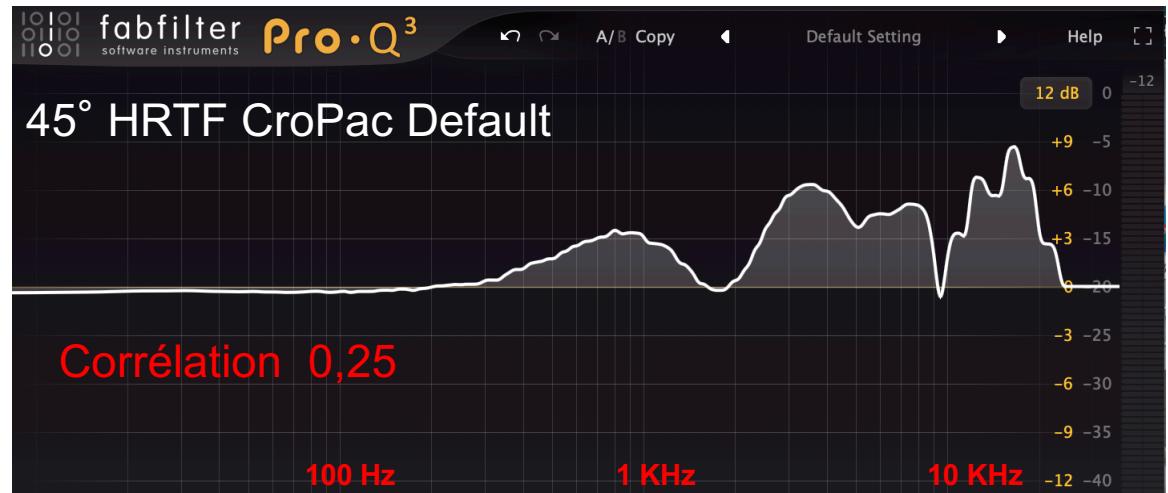
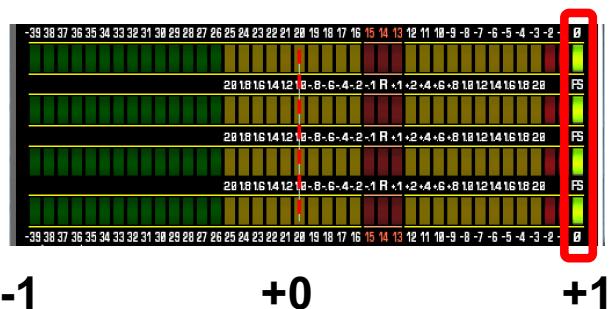


Pour : 45° Azimuth  
0° Elévation

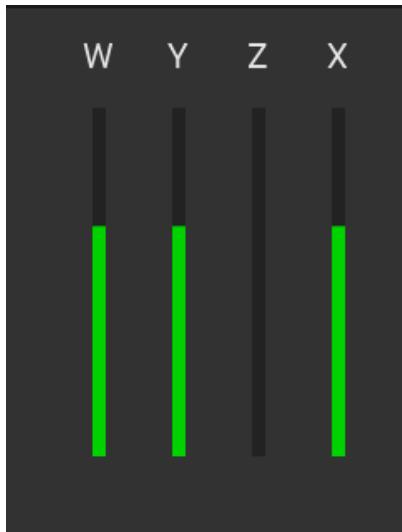


**IN : Format B AmbiX**

Bruit Rose Corrélu sur les 4 canaux  
(Phase à + 1 = mêmes signaux)

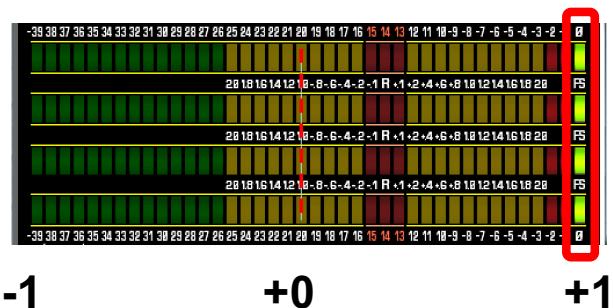


Pour : 45° Azimuth  
0° Elévation

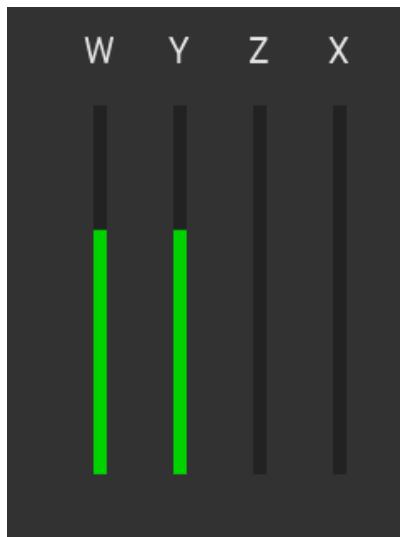


**IN : Format B AmbiX**

Bruit Rose Corrélu sur les 4 canaux  
(Phase à + 1 = mêmes signaux)

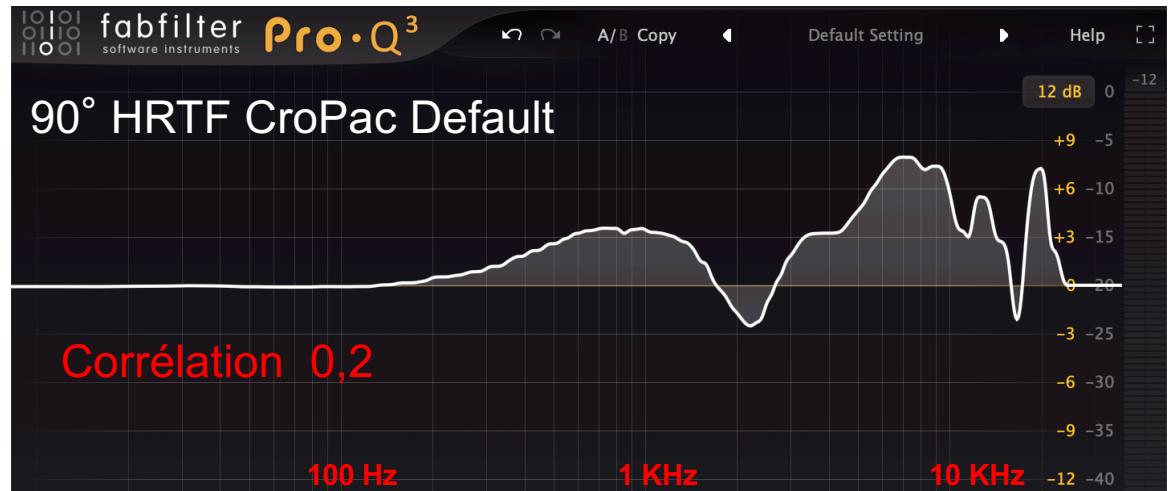
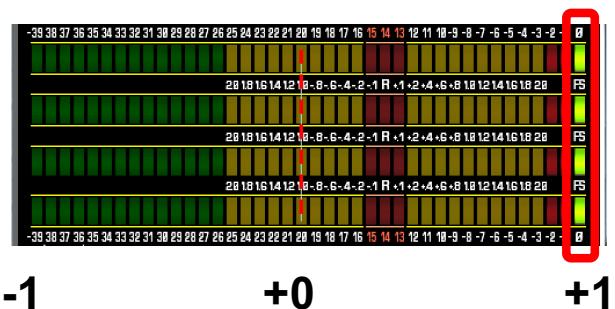


Pour : 90° Azimuth  
0° Elévation

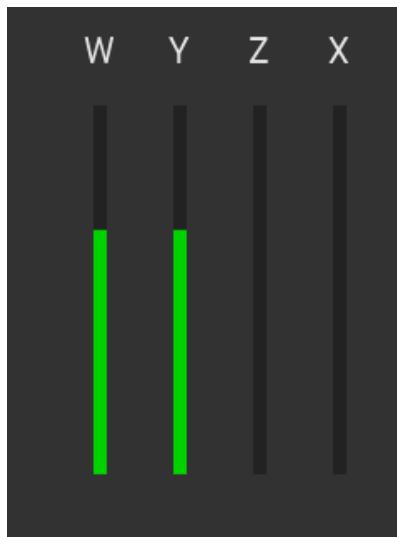


**IN : Format B AmbiX**

Bruit Rose Corrélé sur les 4 canaux  
(Phase à + 1 = mêmes signaux)

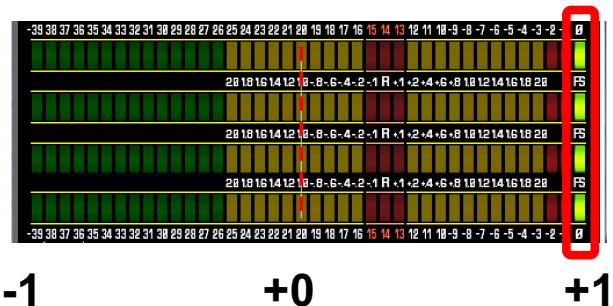


Pour : 90° Azimuth  
0° Elévation



**IN : Format B AmbiX**

Bruit Rose Corrélu sur les 4 canaux  
(Phase à + 1 = mêmes signaux)



# Annexes :



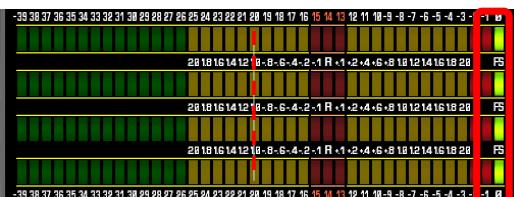
<https://www.dearvr.com/products/ambi-micro>

gratuit

# DEAR VR AMBI MICRO

**IN : Format-B AmbiX**

Bruit Rose Corrélé sur les 4 canaux  
(Phase à + 1 = mêmes signaux)



-1            +0            +1

**0° azimut 0° élévation**





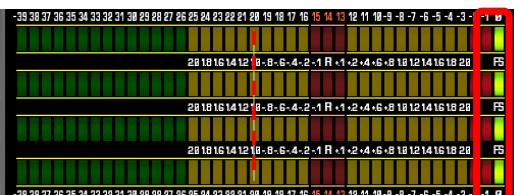
<https://www.dearvr.com/products/ambi-micro>

gratuit

# DEAR VR AMBI MICRO

**IN : Format-B AmbiX**

Bruit Rose Corrélé sur les 4 canaux  
(Phase à + 1 = mêmes signaux)



-1            +0            +1

**45° azimut 45° élévation**



# PARAMETRIC FIRST-ORDER AMBISONIC DECODING FOR HEADPHONES UTILISING THE CROSS-PATTERN COHERENCE ALGORITHM

Leo McCormack and Symeon Delikaris-Manias

Aalto University, Department of Signal Processing and Acoustics, Espoo, Finland

leo.mccormack@aalto.fi

## ABSTRACT

Binaural ambisonics decoding is a means of reproducing a captured or synthesised sound-field, as described by a spherical harmonic representation, over headphones. The majority of ambisonic decoders proposed to date are based on a signal-independent approach; operating via a linear mapping between the input spherical harmonic signals and the output binaural signals. While this approach is computationally efficient, an impractically high input order is often required to deliver a sufficiently accurate rendition of the original spatial cues to the listener. This is especially problematic, as the vast majority of commercially available Ambisonics microphones are first-order, which ultimately results in numerous perceptual deficiencies during reproduction. Therefore, in this paper, a signal-dependent and parametric binaural ambisonic decoder is proposed, which is specifically intended to reproduce first-order input with high perceptual accuracy. The proposed method assumes a sound-field model of one source and one non-isotropic ambient component per narrow-band. It then employs the Cross-Pattern Coherence (CroPaC) post-filter, in order to segregate these components with improved spatial selectivity. Listening test results indicate that the proposed method, when using first-order input, performs similarly to third-order Ambisonics reproduction.

## 1. INTRODUCTION

Reproduction of synthesised or captured sound-fields is an important component in many immersive audio applications, where flexibility, in terms of both content generation and playback setup, is highly favoured. Methods formulated in the spherical harmonic domain (SHD) [1] are often well-suited to this task, as the recording and reproduction operations may be decoupled; with spherical harmonic signals serving as an intermediary. This SHD-based ecosystem for sound-field capture and reproduction is popularly known as Ambisonics [2], where the generation of spherical harmonic signals and the subsequent reproduction of the sound scene that they describe, is referred to as

ambisonic encoding and decoding, respectively. Regarding the latter, currently proposed decoders may be loosely categorised as either non-parametric (signal-independent) or parametric (signal-dependent). Non-parametric binaural reproduction relies on a complex, frequency-dependent, and linear mapping of the input signals to the binaural channels. Whereas parametric methods operate by imposing a set of assumptions regarding the composition of the sound-field and are signal-dependent. Methods that fall within this latter category often rely on the extraction of perceptually meaningful parameters in the time-frequency domain, with the aim of mapping input signals to the binaural channels in a more informed manner [3]. This paper is primarily concerned with parametric reproduction of first-order spherical harmonic input over headphones. A binaural decoder is proposed for this task, which utilises the Cross-Pattern Coherence (CroPaC) [4] spatial post-filter; in order to segregate the sound-field into one source and one non-isotropic ambient component per time-frequency tile during the analysis stage. The method then employs the optimal mixing approach described in [5], to synthesise the output binaural signals.

### 1.1 Non-parametric binaural ambisonics decoding

Binaural ambisonics decoding is conducted via the application of a matrix of filters, which appropriately maps the input spherical harmonic signals to the binaural channels in a linear manner. Therefore, no time-varying distortions are introduced into the output signals. The decoding filters may be derived by approximating the directivity patterns of the listener's head-related transfer functions (HRTF), using the spherical harmonic basis functions, in a least-squares (LS) sense. However, in order to sufficiently approximate and reproduce these complicated directional patterns, a dense grid of HRTF measurements and a high input order is required; often in the range of 15–20th order. For practical reasons, the input order is typically truncated to a much lower order than that of the spatial order of the HRTF measurement grid. This, in turn, results in direction-dependent timbral colourations in the binaural signals. In addition, Ambisonics reproduction is inherently limited by the spatial resolution of the input format. For lower-orders, this has been found to exhibit numerous perceptual deficiencies, including: localisation ambiguity, comb-filtering effects, poor externalisation, and a loss of envelopment [6–10].

Timbral colourations due to input order truncation es-



© Leo McCormack and Symeon Delikaris-Manias. Licensed under a Creative Commons Attribution 4.0 International License (CC BY 4.0). **Attribution:** Leo McCormack and Symeon Delikaris-Manias. "Parametric first-order ambisonic decoding for headphones utilising the Cross-Pattern Coherence algorithm", 1st EAA Spatial Audio Signal Processing Symposium, Paris, France, 2019.

pecially affect high-frequencies, since the high-frequency energy is predominantly concentrated in the higher-order components. This loss of energy may be compensated for via diffuse-field equalisation [11]. However, this loss of high-frequency energy is largely due to the miss-match between the input order and the spatial order of the measurement grid, which is directly proportional to its density. Therefore, rather than applying post-equalisation filters, one may simply reduce the number of points in the HRTF measurement grid, such that its spatial order is more in-line with that of the input order; as suggested in [12]. This approach is often referred to as Spatial Re-sampling (SPR) or virtual loudspeaker decoding. In this case, rather than assigning high-frequency energy to higher-order components and subsequently discarding it, due to order truncation, the energy is instead aliased back into the lower-order components and preserved. However, while this approach improves upon the perceived timbral short-comings of lower-order binaural Ambisonic reproduction, it does not eliminate them, nor does it address the spatial deficiencies of the method.

The localisation ambiguities associated with lower-order binaural Ambisonic reproduction are due to a degradation of the reproduced binaural cues. There are two key causes for this. The first is due to the inherent low input spatial resolution, which leads to erroneously high signal coherence between the output channels; when generated in a linear manner. Whereas the second is due to the LS decoder itself, as it is unable to sufficiently fit the lower-order spherical harmonic patterns to the highly directive HRTF patterns. To address this latter limitation, an alternative method was proposed in [13], which conducts preliminary time-alignment of the Head-related impulse responses (HRIRs) and performs the LS fitting with an additional diffuse-field coherence constraint. The method essentially exploits prior knowledge of the bandwidth in which the inter-aural level differences (ILDs) are the dominant localisation cues; which is above approximately 1.5 kHz, as described by the Haas effect [14]. By discarding the phase information of the HRTFs at frequencies above 1.5 kHz, the LS fitting instead prioritises the delivery of the correct magnitude responses; rather than the phase. Thus it ultimately yields improved ILD cues and diminished inter-aural time difference (ITD) cues; but in a frequency range where ILD cues are more important for localisation. The same principle was also later employed in [15]. However, while these aforementioned approaches yield considerable improvements over traditional decoders, as shown with formal listening tests [13, 16], their performance with first-order input still deviates from that of higher-orders and directly binauralised scenes. This is especially problematic as the vast majority of commercially available Ambisonic microphones and available content are first-order.

## 1.2 Parametric binaural decoding

The inherent perceptual limitations associated with lower-order Ambisonics are primarily as a result of the erroneously high coherence between the output channels. In order to overcome these limitations, signal-dependent and parametric alternatives have been proposed [17–22]. These

methods employ a sound-field model, which lays out a set of assumptions regarding the composition of the sound-field. The methods operate by extracting perceptually meaningful parameters in the time-frequency domain, and often employ dedicated rendering techniques for different components. The two main challenges when designing a parametric method are therefore: 1) identifying a perceptually robust sound-field model, and 2) employing the appropriate signal processing techniques in order to realise the model, with minimal artefacts incurred.

The most well-known and established parametric reproduction method is Directional Audio Coding (DirAC) [17], which employs a sound-field model consisting of one plane-wave and one diffuseness estimate per time-frequency tile. These parameters are derived from the active-intensity vector, in the case of first-order input. The plane-wave components are panned directly to the loudspeakers using Vector-base Amplitude Panning (VBAP) [23], and the diffuse components are sent to all loudspeakers and decorrelated. More recent formulations of DirAC also allow for multiple plane-wave and diffuseness estimates via spatially-localised active-intensity vectors, using higher-order input [18, 19]. In [24], a post-filter was proposed, which adaptively mixes between linearly decoded output and DirAC rendered outputs, in order to improve the output signal fidelity.

High Angular Resolution plane-wave Expansion (HARPEX) [20], is another example of a parametric method, which operates by extracting two plane-wave components per frequency using first-order input. The Sparse-Recovery method [21] extracts a number of plane-waves, which corresponds to up to half the number of input channels of arbitrary order. The COding and Multi-Parameterisation of Ambisonic Sound Scenes (COMPASS) method [25] also extracts multiple source components; up to half the number of input channels. However, it employs an additional residual stream that encapsulates the remaining diffuse and ambient components in the scene. An alternative parameterisation of the sound-field was also presented in [26, 27], which circumvents the modelling of the sound-field with conventional parameters, such as source direction or diffuseness. It considers only the perceived quality of the individual output channels and the perceived quality of the spatial attributes of the reproduction. In [22] an adaptive LS-binaural decoder was proposed, which constrains the covariance matrix of the decoded audio to reflect that of the covariance matrix of the listeners HRTFs. The method does not require explicit estimation of the direct/diffuse balance, nor does it need to employ de-correlation.

### 1.2.1 Parametric binaural decoding with optimal mixing

The main point of criticism regarding parametric reproduction methods is the incursion of time-varying artefacts. These occur either due to the input scene not conforming to the assumed sound-field model, or due to the rendering techniques not being sufficiently robust. Therefore, in [28], an *optimal mixing* technique was proposed, which attempts to synthesise the output using a linear combination of linearly decoded prototype signals, as much as possible; thus

retaining much of the high single-channel fidelity in the output. The method then employs decorrelation only if the output inter-channel dependencies still deviate from the target, thus mitigating potential decorrelation artefacts; such as the temporal smearing of transients.

The approach is formulated in the covariance domain and relies on the construction of time-varying narrow-band target covariance matrices, which are dictated by the sound-field model of the parametric method. It then employs traditional linearly decoded audio as a prototype. The approach has been employed previously in [19] using the DirAC sound-field model, and also in the closed-form solution of [22] for binaural reproduction. However, in principle, any sound-field model may employ this optimal mixing approach for the synthesis stage.

### 1.3 Motivation for this work

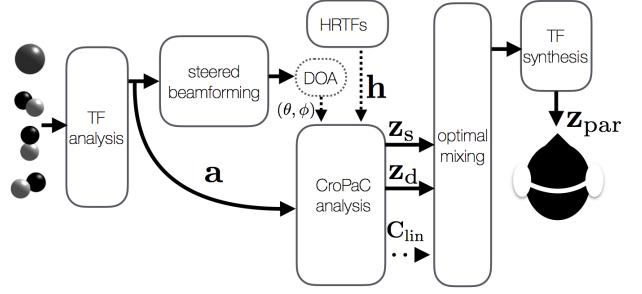
Further development of first-order ambisonic decoders is of particular interest; given the prevalence of first-order commercially available Ambisonic microphone arrays and material accessible on the internet. Recent advancements have shown improvements in the perceived performance of linear binaural ambisonic decoding [13, 15]. However, these decoders still rely on input orders which may be considered impractical for wide-spread adoption today. Especially given the quadratic growth in the number of channels (and microphones) with increasing order, leading to more expensive arrays and higher bandwidths. Therefore, the need for robust signal-dependent parametric alternatives for lower-orders is well established.

In this paper, a first-order binaural decoder is proposed<sup>1</sup>, which employs a parametric sound-field model that assumes one source and one directional ambient component per time-frequency tile. The model is inspired by the multi-source and non-isotropic ambient model employed by the COMPASS method [25]. However, given the low-resolution of first-order input, the proposed method employs an additional CroPaC spatial post-filter [4] to better isolate the source component, and an inverse CroPaC post-filter to obtain the ambient component. The proposed method also extracts instantaneous direction-of-arrival (DoA) estimates via steered-response power (SRP) activity-maps and peak-finding; thus forgoing the need for long temporal averaging of input covariance matrices, as required by sub-space alternatives. The method also employs the optimal mixing solution of [5], in order to minimise the amount of decorrelated signal energy in the output, and synthesise the target inter-channel dependencies in a linear manner, as much as possible.

### 1.4 Organisation of the paper

Section 2 describes linear binaural ambisonic decoding, which is employed as a basis for the proposed method. Section 3 details the proposed method. Subjective evaluation of the proposed algorithm, through listening tests, is described in Section 4, and Section 5 concludes the paper.

<sup>1</sup> A VST audio plug-in of the proposed decoder may be found here: <http://research.spa.aalto.fi/publications/papers/sasp19-parametric/>



**Figure 1:** Overall block diagram of the proposed method.

## 2. LINEAR BINAURAL AMBISONICS DECODING

Due to the signal- and frequency-dependent nature of the proposed algorithm (described in Section 3) it is assumed that the input ambisonic signals have been transformed into the time-frequency ( $t, f$ ) domain; where  $t$  and  $f$  denote the time and frequency indices, respectively. The ambisonic signals,  $\mathbf{a}$ , may be synthesised by mapping monophonic signals onto spherical harmonic basis functions or captured utilising a microphone array with subsequent suitable encoding

$$\mathbf{a} = [a_{00}, a_{1(-1)}, a_{10}, \dots, a_{N(N-1)}, a_{NN}]^T \in \mathbb{C}^{(N+1)^2 \times 1}, \quad (1)$$

where  $a_{nm}$  are the individual ambisonic signals for each order,  $n$ , and degree,  $m$ , up to the maximum order,  $N$ . It is assumed, henceforth, that the input ambisonic signals conform to the ortho-normalised (N3D) and Ambisonic Channel Numbering (ACN) conventions.

Ambisonic signals may be decoded for headphone playback as

$$\mathbf{z}_{\text{lin}}(t, f) = \mathbf{D}_{\text{bin}}(f)\mathbf{a}(t, f), \quad (2)$$

where  $\mathbf{z}_{\text{lin}}(t, f) \in \mathbb{C}^{2 \times 1}$  are the output binaural signals, and  $\mathbf{D}_{\text{bin}}(f) \in \mathbb{C}^{2 \times (N+1)^2}$  is a binaural ambisonic decoding matrix, derived using one of a number of approaches [12, 13, 15].

## 3. PROPOSED PARAMETRIC BINAURAL AMBISONICS DECODER

The proposed first-order decoder employs a sound-field model comprising one source component and one non-isotropic ambient component per time-frequency tile. The method first estimates the source DoA via steered-response power (SRP) beamforming and subsequent peak-finding. The source stream is then segregated by steering a beamformer toward the estimated DoA, and employing an additional CroPaC post-filtering operation to improve its spatial selectivity. The ambient stream is then simply the residual, once the source component has been subtracted from the input sound-field. The two streams are then binauralised and fed into an optimal mixing unit, along with the ambisonic prototype covariance matrix, in order to generate the binaural output. A block diagram of the proposed method is depicted in Fig. 1.

### 3.1 Analysis

#### 3.1.1 DoA estimation

A DoA estimator based on SRP beamforming and peak-finding [29] is suggested for the proposed rendering method. Note that the chosen scanning grid should preferably take the angular resolution of human hearing into account [14]. This approach yields instantaneous estimates of the source direction,  $(\theta_s, \phi_s)$ , which is in keeping with the instantaneous CroPaC post-filter values. Therefore, the need for temporal averaging of input covariance matrices is avoided in the analysis stage, as would be required by subspace-based alternatives.

#### 3.1.2 CroPaC post-filter

The cross-correlation between the omni and dipole may be utilised as a spatial post-filter [4]

$$G(t, f) = \frac{2}{\sqrt{3}} \Re[a_{00}^*(t, f)a_{11}(t, f)], \quad (3)$$

where  $\Re$  denotes the real operator and  $*$  denotes the complex conjugate. This value is then normalised with the input sound-field energy and half-wave rectified

$$\hat{G}(t, f) = \max \left[ \frac{G(t, f)}{|a_{00}(t, f)|^2 + \sum_{-1}^1 |a_{1m}(t, f)|/\sqrt{3}|^2}, \lambda \right], \quad (4)$$

where  $\lambda \in [0, \dots, 1]$  is a parameter which influences the severity of the post-filter, similarly to the spectral floor of a traditional post-filter. Note that the spatial selectivity of the CroPaC post-filter is sharper than a conventional first-order beam pattern, due to the multiplication (rather than summation) of the two signals in (3).

The CroPaC spatial filter may be steered in the source direction, by first rotating the spherical harmonic signals using an appropriate rotation matrix,  $\mathbf{M}_{\text{rot}}(\theta_s, \phi_s) \in \mathbb{R}^{(N+1)^2 \times (N+1)^2}$ , as

$$\mathbf{a}_{\text{rot}}(t, f) = \mathbf{M}_{\text{rot}}(\theta_s, \phi_s)\mathbf{a}(t, f), \quad (5)$$

where  $\mathbf{a}_{\text{rot}}(t, f) \in \mathbb{C}^{(N+1)^2 \times 1}$  are the resulting rotated signals, which are then employed for the post-filter estimation (3). For details regarding the calculation of this rotation matrix, the reader is directed to [30].

### 3.2 Synthesis

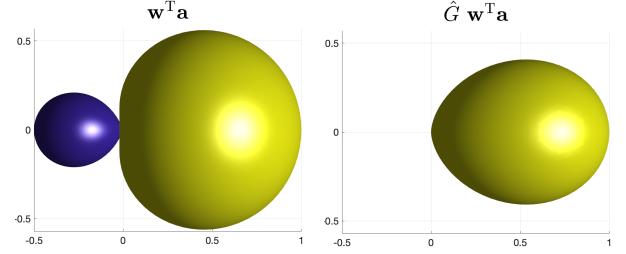
The source stream,  $\mathbf{z}_s(t, f)$ , is obtained by directly applying the HRTF gains to the extracted source signal as

$$\mathbf{z}_s(t, f) = \mathbf{h} \frac{\hat{G}(t, f)}{(N+1)^2} \mathbf{w}^T \mathbf{a}(t, f), \quad (6)$$

where  $\mathbf{h} \in \mathbb{C}^{2 \times 1}$  and  $\mathbf{w} \in \mathbb{R}^{(N+1)^2 \times 1}$  are HRTF gains and static beamforming weights [1], corresponding to the analysed DoA at each time-frequency tile, respectively. A visual depiction of the improved spatial selectivity of the source beamformer with the post-filter is given in Fig. 2.

The ambient stream,  $\mathbf{z}_d(t, f)$ , is then obtained by subtracting the source signal from the input scene, and decoding it to headphones using a binaural ambisonic decoder

$$\mathbf{z}_d(t, f) = \mathbf{D}_{\text{bin}}[\mathbf{a}(t, f) - \frac{\hat{G}(t, f)}{(N+1)^2} \mathbf{y} \mathbf{w}^T \mathbf{a}(t, f)], \quad (7)$$



**Figure 2:** A first-order hyper cardioid beamformer without (left) and with (right) the CroPaC post-filter ( $\lambda = 0$ ).

where  $\mathbf{y} \in \mathbb{R}^{(N+1)^2 \times 1}$  are the spherical harmonic weights for the analysed DoA.

Optionally, the ambient stream may be decorrelated, in order to further minimise the inter-channel coherence between the binaural channels

$$\tilde{\mathbf{z}}_d(t, f) = \mathcal{D}[\mathbf{z}_d(t, f)] \quad (8)$$

where  $\mathcal{D}[\cdot]$  denotes a decorrelation operation on the enclosed signals.

#### 3.2.1 Optimal mixing

Rather than summing the source (6) and ambient (8) streams together, to acquire the binaural output, an alternative synthesis approach is suggested. This approach is based on the covariance domain framework, termed here as *optimal mixing*, and is described in [5, 28]. The method employs linearly decoded signals,  $\mathbf{z}_{\text{lin}}(t, f)$ , as a prototype, which has the following base-line covariance matrix (note that the time-frequency indices have been omitted for the brevity of notation)

$$\mathbf{C}_{\text{lin}} = \mathbb{E}[\mathbf{z}_{\text{lin}} \mathbf{z}_{\text{lin}}^H] = \mathbf{D}_{\text{bin}} \mathbf{C}_a \mathbf{D}_{\text{bin}}^H, \quad (9)$$

where  $\mathbf{C}_a = \mathbb{E}[\mathbf{a} \mathbf{a}^H] \in \mathbb{C}^{(N+1)^2 \times (N+1)^2}$  is the covariance matrix of the input signals.

A time-varying and narrow-band target covariance matrix is then required, which may be defined in this case as

$$\mathbf{C}_{\text{target}} = \mathbb{E}[(\mathbf{z}_s + \mathbf{z}_d)(\mathbf{z}_s + \mathbf{z}_d)^H]. \quad (10)$$

The optimal mixing solution then provides the values for matrices  $\mathbf{A} \in \mathbb{C}^{2 \times 2}$  and  $\mathbf{B} \in \mathbb{C}^{2 \times 2}$ , in the following equation

$$\mathbf{A} \mathbf{C}_{\text{lin}} \mathbf{A}^H + \mathbf{B} \tilde{\mathbf{C}}_{\text{lin}} \mathbf{B}^H = \mathbf{C}_{\text{target}} \quad (11)$$

where  $\tilde{\mathbf{C}}_{\text{lin}} = \text{diag}[\mathcal{D}[\mathbf{z}_{\text{lin}}] \mathcal{D}[\mathbf{z}_{\text{lin}}]^H]$  is a diagonal matrix consisting of the diagonal entries of the covariance matrix of a decorrelated version of the linearly decoded prototype. More information regarding the derivation and calculation of these mixing matrices is given in [5].

The output audio  $\mathbf{z}_{\text{par}}$  may then be obtained as

$$\mathbf{z}_{\text{par}} = \mathbf{A} \mathbf{z}_{\text{lin}} + \mathbf{B} \mathcal{D}[\mathbf{z}_{\text{lin}}], \quad (12)$$

which ideally exhibits all of the target inter-channel dependencies, as dictated by the target covariance matrix  $\mathbb{E}[\mathbf{z}_{\text{par}} \mathbf{z}_{\text{par}}^H] \simeq \mathbf{C}_{\text{target}}$ . However, note that in practice, due to the need for regularisation of the input covariance matrices, the solution is never exact.

#### 4. EVALUATION

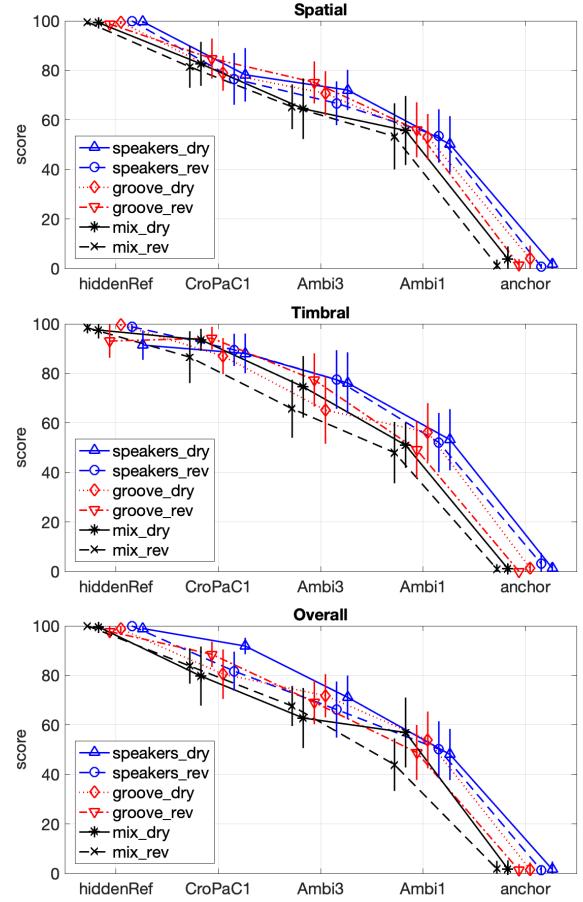
A multiple-stimulus test was conducted, in order to compare the perceived performance of the proposed approach (*CroPaC1*), with both first-order (*Ambi1*) and third-order (*Ambi3*) spatial re-sampling ambisonic decoders [12]. Note that the optimal mixing approach (12) was employed, using the same ambisonic decoder for the prototype  $\mathbf{z}_{\text{lin}}$ .

Synthetic test scenes were created as follows: multiple speakers (*speakers*), a modern funk band (*groove*), and a mix comprising of a speaker, clapping, water fountain and piano (*mix*). The individual sources were binauralised directly for the reference test cases, and encoded into first- and third- order signals and passed through their respective decoders for the *\_dry* test cases. A shoebox image-source room simulator<sup>2</sup> was employed to obtain reverberant (*\_rev*) counterparts. The room simulator was configured to resemble a small hall. The image-sources arriving at the receiver position were binauralised directly for the reference cases, and encoded into spherical harmonic signals for the reverberant test cases.

The evaluation consisted of three parts, which addressed the spatial, timbral, and overall differences between the methods. For the evaluation of **spatial** differences, the mean spectra of the reference was imposed onto the test cases. Therefore, timbral differences between the methods were greatly mitigated, but still retained their original spatial characteristics. The test subjects were explicitly instructed to ignore any remaining timbral differences, which may not have been addressed by the equalisation. The anchor was obtained as an omni-directional spherical harmonic component, replicated to each binaural channel and also equalised. For evaluating the **timbral** differences, the mean spectra of each test case, was imposed onto an omni-directional signal and sent to both left and right channels. Therefore, the spatial differences between the methods were eliminated, thus retaining only the timbral differences. The anchor was obtained as the mean spectra of the reference, replicated to each binaural channel and low-passed filtered at 4 kHz. Finally, for assessing the **overall** differences, only the broad-band RMS of the reference was used to normalise the test cases. Therefore, all timbral and spatial differences between test cases remained, and the test participants graded the samples based on their subjective weighting of the reproduced attributes.

A total of 14 expert listeners participated in the listening tests in purpose-built headphone booths. The present authors did not take part in the tests. The means and 95% confidence intervals of the results are shown in Fig. 3. It can be observed that the spatial and timbral characteristics of the reference were more closely reproduced using the proposed method, when compared to conventional ambisonics with the same first-order input. Furthermore, the proposed method yielded scores more inline with that of third-order ambisonics. It should be highlighted that third-order ambisonics employs four times the number of input channels than that of the proposed method. This, therefore, represents a significant reduction in bandwidth, without compromising the perceived spatial accuracy or fidelity.

<sup>2</sup> The shoe-box room simulator employed for the reverberant test cases may be found here: <https://github.com/polarch/shoebox-roomsim>



**Figure 3:** The means and 95% confidence intervals for each individual sound scene. The evaluation criteria were: spatial only, timbral only, and overall (top-bottom).

#### 5. CONCLUSION

This paper has proposed a first-order parametric binaural ambisonic decoder, which employs a sound-field model comprising one source and one directional ambient component per time-frequency tile. The proposed approach first isolates the source components using a spherical harmonic domain beamformer, modulated by the Cross-Pattern Coherence (*CroPaC*) spatial post-filter. The ambient stream is then simply the residual, once the source components have been subtracted from the input sound-field. The proposed approach is inspired by the COMPASS method [25]. However, along with the *CroPaC* post-filter, it also employs instantaneous source direction estimation and synthesises the output in a linear manner as much as possible; in order to improve the fidelity of the output signals. Formal listening tests indicate that the proposed first-order decoder performs similarly to (or exceeds) third-order spatial re-sampling ambisonics decoding, in terms of the perceived spatial and timbral attributes of the reproduction.

#### 6. ACKNOWLEDGEMENTS

This research has received funding from the Aalto ELEC Doctoral School. Many thanks are also extended to Dr. Archontis Politis for the insightful discussions regarding the method.

## 7. REFERENCES

- [1] B. Rafaely, *Fundamentals of spherical array processing*, vol. 8. Springer, 2015.
- [2] M. A. Gerzon, "Periphony: With-height sound reproduction," *Journal of the Audio Engineering Society*, vol. 21, no. 1, pp. 2–10, 1973.
- [3] V. Pulkki, S. Delikaris-Manias, and A. Politis, "Parametric time-frequency domain spatial audio," 2018.
- [4] S. Delikaris-Manias and V. Pulkki, "Cross pattern coherence algorithm for spatial filtering applications utilizing microphone arrays," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 21, no. 11, pp. 2356–2367, 2013.
- [5] J. Vilkamo and T. Backstrom, "Time–frequency processing: Methods and tools," in *Parametric Time-Frequency Domain Spatial Audio*, pp. 3–23, John Wiley & Sons, 2018.
- [6] O. Santala, H. Vertanen, J. Pekonen, J. Oksanen, and V. Pulkki, "Effect of listening room on audio quality in ambisonics reproduction," in *Audio Engineering Society Convention I26*, Audio Engineering Society, 2009.
- [7] S. Braun and M. Frank, "Localization of 3d ambisonic recordings and ambisonic virtual sources," in *1st International Conference on Spatial Audio,(Detmold)*, 2011.
- [8] A. Avni, J. Ahrens, M. Geier, S. Spors, H. Wierstorf, and B. Rafaely, "Spatial perception of sound fields recorded by spherical microphone arrays with varying spatial resolution," *The Journal of the Acoustical Society of America*, vol. 133, no. 5, pp. 2711–2721, 2013.
- [9] S. Bertet, J. Daniel, E. Parizet, and O. Warusfel, "Investigation on localisation accuracy for first and higher order ambisonics reproduced sound sources," *Acta Acustica united with Acustica*, vol. 99, no. 4, pp. 642–657, 2013.
- [10] P. Stitt, S. Bertet, and M. van Walstijn, "Off-centre localisation performance of ambisonics and hoa for large and small loudspeaker array radii," *Acta Acustica united with Acustica*, vol. 100, no. 5, pp. 937–944, 2014.
- [11] Z. Ben-Hur, F. Brinkmann, J. Sheaffer, S. Weinzierl, and B. Rafaely, "Spectral equalization in binaural signals represented by order-truncated spherical harmonics," *The Journal of the Acoustical Society of America*, vol. 141, no. 6, pp. 4087–4096, 2017.
- [12] B. Bernschütz, A. V. Giner, C. Pörschmann, and J. Arend, "Binaural reproduction of plane waves with reduced modal order," *Acta Acustica united with Acustica*, vol. 100, no. 5, pp. 972–983, 2014.
- [13] M. Zaunschirm, C. Schörkhuber, and R. Höldrich, "Binaural rendering of ambisonic signals by head-related impulse response time alignment and a diffuseness constraint," *The Journal of the Acoustical Society of America*, vol. 143, no. 6, pp. 3616–3627, 2018.
- [14] J. Blauert, *Spatial hearing: the psychophysics of human sound localization*. MIT press, 1997.
- [15] C. Schörkhuber, M. Zaunschirm, and R. Höldrich, "Binaural rendering of ambisonic signals via magnitude least squares," in *Proceedings of the DAGA*, vol. 44, pp. 339–342, 2018.
- [16] H. Lee, M. Frank, and F. Zotter, "Spatial and timbral fidelities of binaural ambisonics decoders for main microphone array recordings," in *Audio Engineering Society Conference: 2019 AES International Conference on Immersive and Interactive Audio*, Audio Engineering Society, 2019.
- [17] V. Pulkki, "Directional audio coding in spatial sound reproduction and stereo upmixing," in *Audio Engineering Society Conference: 28th International Conference: The Future of Audio Technology—Surround and Beyond*, Audio Engineering Society, 2006.
- [18] A. Politis, J. Vilkamo, and V. Pulkki, "Sector-based parametric sound field reproduction in the spherical harmonic domain," *IEEE Journal of Selected Topics in Signal Processing*, vol. 9, no. 5, pp. 852–866, 2015.
- [19] A. Politis, L. McCormack, and V. Pulkki, "Enhancement of ambisonic binaural reproduction using directional audio coding with optimal adaptive mixing," in *IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*, 2017.
- [20] S. Berge and N. Barrett, "High angular resolution planewave expansion," in *Proc. of the 2nd International Symposium on Ambisonics and Spherical Acoustics May*, pp. 6–7, 2010.
- [21] A. Wabnitz, N. Epain, and C. T. Jin, "A frequency-domain algorithm to upscale ambisonic sound scenes," in *2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 385–388, IEEE, 2012.
- [22] C. Schörkhuber and R. Höldrich, "Linearly and quadratically constrained least-squares decoder for signal-dependent binaural rendering of ambisonic signals," in *Audio Engineering Society Conference: 2019 AES International Conference on Immersive and Interactive Audio*, Audio Engineering Society, 2019.
- [23] V. Pulkki, "Virtual sound source positioning using vector base amplitude panning," *Journal of the audio engineering society*, vol. 45, no. 6, pp. 456–466, 1997.
- [24] O. Thiergart, G. Milano, and E. A. Habets, "Combining linear spatial filtering and non-linear parametric processing for high-quality spatial sound capturing," in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 571–575, IEEE, 2019.
- [25] A. Politis, S. Tervo, and V. Pulkki, "Compass: Coding and multidirectional parameterization of ambisonic sound scenes," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 6802–6806, IEEE, 2018.
- [26] J. Vilkamo and S. Delikaris-Manias, "Perceptual reproduction of spatial sound using loudspeaker-signal-domain parametrization," *IEEE/ACM Transactions on Audio, Speech and Language Processing (TASLP)*, vol. 23, no. 10, pp. 1660–1669, 2015.
- [27] S. Delikaris-Manias, J. Vilkamo, and V. Pulkki, "Parametric binaural rendering utilizing compact microphone arrays," in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 629–633, IEEE, 2015.
- [28] J. Vilkamo, T. Bäckström, and A. Kuntz, "Optimized covariance domain framework for time–frequency processing of spatial audio," *Journal of the Audio Engineering Society*, vol. 61, no. 6, pp. 403–411, 2013.
- [29] D. P. Jarrett, E. A. Habets, and P. A. Naylor, "3d source localization in the spherical harmonic domain using a pseudointensity vector," in *2010 18th European Signal Processing Conference*, pp. 442–446, IEEE, 2010.
- [30] J. Ivanic and K. Ruedenberg, "Rotation matrices for real spherical harmonics. direct determination by recursion," *The Journal of Physical Chemistry A*, vol. 102, no. 45, pp. 9099–9100, 1998.

# TRADUCTION FR

## DÉCODAGE AMBISONIQUE PARAMÉTRIQUE DE PREMIER ORDRE POUR CASQUE D'ÉCOUTE UTILISANT L'ALGORITHME DE COHÉRENCE CROSS-PATTERN

**Leo McCormack et Symeon Delikaris-Manias**

Université Aalto, Département de traitement du signal et d'acoustique, Espoo, Finlande

[leo.mccormack@aalto.fi](mailto:leo.mccormack@aalto.fi)

### ABSTRAIT

Le décodage binaural ambisonique est un moyen de reproduire un champ sonore capturé ou synthétisé, tel que décrit par une représentation harmonique sphérique, au casque. La majorité des décodeurs ambisoniques proposés à ce jour sont basés sur une approche indépendante du signal ; fonctionnant via un mappage linéaire entre les signaux harmoniques sphériques d'entrée et les signaux binauraux de sortie. Bien que cette approche soit efficace sur le plan informatique, un ordre d'entrée trop élevé est souvent nécessaire pour fournir à l'auditeur un rendu suffisamment précis des repères spatiaux d'origine. Ceci est particulièrement problématique, car la grande majorité des microphones Ambisonics disponibles dans le commerce sont de premier ordre, ce qui entraîne finalement de nombreuses déficiences perceptuelles lors de la reproduction. Par conséquent, dans cet article, un décodeur ambisonique binaural dépendant du signal et paramétrique est proposé, qui est spécifiquement destiné à reproduire une entrée de premier ordre avec une précision perceptuelle élevée. La méthode proposée suppose un modèle de champ sonore d'une source et d'une composante ambiante non isotrope par bande étroite. Il utilise ensuite le post-filtre Cross-Pattern Coherence (CroPaC), afin de séparer ces composants avec une meilleure sélectivité spatiale. Les résultats des tests d'écoute indiquent que la méthode proposée, lorsqu'elle utilise une entrée de premier ordre, fonctionne de manière similaire à la reproduction Ambisonics de troisième ordre. Il utilise ensuite le post-filtre Cross-Pattern Coherence (CroPaC), afin de séparer ces composants avec une meilleure sélectivité spatiale. Les résultats des tests d'écoute indiquent que la méthode proposée, lorsqu'elle utilise une entrée de premier ordre, fonctionne de manière similaire à la reproduction Ambisonics de troisième ordre.

## INTRODUCTION

La reproduction de champs sonores synthétisés ou capturés est un élément important dans de nombreuses applications audio immersives, où la flexibilité, tant en termes de génération de contenu que de configuration de lecture, est très favorisée. Les méthodes formulées dans le domaine harmonique sphérique (SHD) [1] sont souvent bien adaptées à cette tâche, car les opérations d'enregistrement et de reproduction peuvent être découplées ; avec des signaux harmoniques sphériques servant d'intermédiaire. Cet écosystème basé sur le SHD pour la capture et la reproduction du champ sonore est communément appelé Ambisonics [2], où la génération de signaux harmoniques sphériques et la reproduction ultérieure de la scène sonore qu'ils décrivent, est appelée codage et décodage ambisonique, respectivement. En ce qui concerne ce dernier, les décodeurs actuellement proposés peuvent être grossièrement classés comme non paramétriques (indépendants du signal) ou paramétriques (dépendants du signal). La reproduction binaurale non paramétrique repose sur un mappage complexe, dépendant de la fréquence et linéaire des signaux d'entrée vers les canaux binauraux. Alors que les méthodes paramétriques fonctionnent en imposant un ensemble d'hypothèses concernant la composition du champ sonore et dépendent du signal. Les méthodes qui entrent dans cette dernière catégorie reposent souvent sur l'extraction de paramètres perceptuellement significatifs dans le domaine temps-fréquence, dans le but de mapper les signaux d'entrée sur les canaux binauraux de manière plus informée [3]. Cet article s'intéresse principalement à la reproduction paramétrique de l'entrée harmonique sphérique du premier ordre sur un casque. Un décodeur binaural est proposé pour cette tâche, qui utilise le post-filtre spatial Cross-Pattern Coherence (CroPaC) [4] ; afin de séparer le champ sonore en une source et une composante ambiante non isotrope par pavé temps-fréquence lors de la phase d'analyse. La méthode utilise ensuite l'approche de mélange optimal décrite dans [5], pour synthétiser les signaux binauraux de sortie.

## Décodage ambisonique binaural non paramétrique

Le décodage binaural ambisonique est effectué via l'application d'une matrice de filtres, qui mappe de manière appropriée les signaux harmoniques sphériques d'entrée sur les canaux binauraux de manière linéaire. Par conséquent, aucune distorsion variant dans le temps n'est introduite dans les signaux de sortie. Les filtres de décodage peuvent être dérivés en approximant les modèles de directivité des fonctions de transfert liées à la tête de l'auditeur (HRTF), en utilisant les fonctions de base harmonique sphérique, dans le sens des moindres carrés (LS). Cependant, afin d'approximer et de reproduire suffisamment ces modèles directionnels compliqués, une grille dense de mesures HRTF et un ordre d'entrée élevé sont nécessaires ; souvent dans la gamme de l'ordre 15-20ème. Pour des raisons pratiques, l'ordre d'entrée est typiquement tronqué à un ordre bien inférieur à celui de l'ordre spatial de la grille de mesure HRTF. Cette, à son tour, entraîne des colorations timbrales dépendantes de la direction dans les signaux binauraux. De plus, la reproduction Ambisonics est intrinsèquement limitée par la résolution spatiale du format d'entrée. Pour les ordres inférieurs, il a été constaté que cela présentait de nombreuses déficiences perceptives, notamment : une ambiguïté de localisation, des effets de filtrage en peigne, une mauvaise externalisation et une perte d'enveloppement [6-10].

Colorations de timbre dues à la troncature de l'ordre d'entrée affectent particulièrement les hautes fréquences, puisque l'énergie haute fréquence est principalement concentrée dans les composants d'ordre supérieur. Cette perte d'énergie peut être compensée par une égalisation en champ diffus [11]. Cependant, cette perte d'énergie haute fréquence est en grande partie due au décalage entre l'ordre d'entrée et l'ordre spatial de la grille de mesure, qui est directement proportionnel à sa densité. Par conséquent, plutôt que d'appliquer des filtres de post-égalisation, on peut simplement réduire le nombre de points dans la grille de mesure HRTF, de sorte que son ordre spatial soit plus conforme à celui de l'ordre d'entrée ; comme suggéré dans [12]. Cette approche est souvent appelée rééchantillonnage spatial (SPR) ou décodage de haut-parleur virtuel. Dans ce cas, plutôt que d'attribuer de l'énergie haute fréquence à des composants d'ordre supérieur et de la rejeter par la suite, en raison de la troncature de l'ordre, l'énergie est à la place réintroduite dans les composants d'ordre inférieur et préservée. Cependant, bien que cette approche améliore les défauts timbraux perçus de la reproduction ambisonique binaurale d'ordre inférieur, elle ne les élimine pas, ni ne résout les défauts spatiaux de la méthode.

Les ambiguïtés de localisation associées à la reproduction ambisonique binaurale d'ordre inférieur sont dues à une dégradation des signaux binauraux reproduits. Il y a deux causes principales à cela. Le premier est dû à la faible résolution spatiale d'entrée inhérente, qui conduit à une cohérence de signal faussement élevée entre les canaux de sortie ; lorsqu'il est généré de manière linéaire. Alors que la seconde est due au décodeur LS lui-même, car il est incapable d'adapter suffisamment les modèles harmoniques sphériques d'ordre inférieur aux modèles HRTF hautement directifs. Pour remédier à cette dernière limitation, une méthode alternative a été proposée dans [13], qui effectue un alignement temporel préliminaire des réponses impulsionales liées à la tête (HRIR) et effectue l'ajustement LS avec une contrainte de cohérence de champ diffus supplémentaire. La méthode exploite essentiellement la connaissance a priori de la bande passante dans laquelle les différences de niveaux inter- auriculaires (ILD) sont les indices de localisation dominants ; qui est supérieure à environ 1,5 kHz, comme décrit par l'effet Haas [14]. En rejetant les informations de phase des HRTF à des fréquences supérieures à 1,5 kHz, le montage LS donne plutôt la priorité à la fourniture des réponses d'amplitude correctes ; plutôt que la phase. Ainsi, il produit finalement des signaux ILD améliorés et des signaux de différence de temps interauraux (ITD) diminués ; mais dans une gamme de fréquences où les indices ILD sont plus importants pour la localisation. Le même principe a également été utilisé plus tard dans [15]. Cependant, bien que ces approches susmentionnées apportent des améliorations considérables par rapport aux décodeurs traditionnels, comme le montrent les tests d'écoute formels [13, 16], leur performance avec une entrée de premier ordre s'écarte encore de celle des scènes d'ordre supérieur et directement binaurisées. Ceci est particulièrement problématique car la grande majorité des microphones ambisoniques disponibles dans le commerce et du contenu disponible sont de premier ordre.

## Décodage binaural paramétrique

Les limitations perceptuelles inhérentes associées à l'ambisonique d'ordre inférieur sont principalement dues à la cohérence erronément élevée entre les canaux de sortie. Afin de surmonter ces limitations, des alternatives dépendantes du signal et paramétriques ont été proposées [17–22]. Ces méthodes utilisent un modèle de champ sonore, qui établit un ensemble d'hypothèses concernant la composition du champ sonore. Les procédés fonctionnent en extrayant des paramètres perceptuellement significatifs dans le domaine temps- fréquence et utilisent souvent des techniques de rendu dédiées pour différents composants. Les deux principaux défis lors de la conception d'une méthode paramétrique sont donc : 1) l'identification d'un modèle de champ sonore perceptiblement robuste, et 2) l'utilisation des techniques de traitement du signal appropriées afin de réaliser le modèle, avec un minimum d'artefacts encourus.

La méthode de reproduction paramétrique la plus connue et la plus établie est le codage audio directionnel (DirAC) [17], qui utilise un modèle de champ sonore composé d'une onde plane et d'une estimation de la diffusion par pavé temps-fréquence. Ces paramètres sont dérivés du vecteur d'intensité active, dans le cas d'une entrée de premier ordre. Les composantes d'onde plane sont panoramisées directement vers les haut- parleurs à l'aide du Vector-Base Amplitude Panning (VBAP) [23], et les composantes diffuses sont envoyées à tous les haut- parleurs et décorrélées. Des formulations plus récentes de DirAC permettent également de multiples estimations d'ondes planes et de diffusion via des vecteurs d'intensité active localisés dans l'espace, en utilisant une entrée d'ordre supérieur [18, 19].

Dans [24], un post-filtre a été proposé, qui mélange de manière adaptative entre la sortie décodée linéairement et les sorties rendues DirAC, L'expansion d'onde plane à haute résolution angulaire (HARPEX) [20] est un autre exemple de méthode paramétrique, qui fonctionne en extrayant deux composantes d'onde plane par fréquence en utilisant une entrée de premier ordre. La méthode Sparse-Recovery [21] extrait un certain nombre d'ondes planes, ce qui correspond à jusqu'à la moitié du nombre de canaux d'entrée d'ordre arbitraire.

La méthode COding and Multi- Parameterisation of Ambisonic Sound Scenes (COMPASS) [25] extrait également plusieurs composantes sources ; jusqu'à la moitié du nombre de canaux d'entrée. Cependant, il utilise un flux résiduel supplémentaire qui encapsule les composants diffus et ambients restants dans la scène. Une paramétrisation alternative du champ sonore a également été présentée dans [26, 27], qui contourne la modélisation du champ sonore avec des paramètres conventionnels, tels que la direction de la source ou la diffusion. Il considère uniquement la qualité perçue des canaux de sortie individuels et la qualité perçue des attributs spatiaux de la reproduction. Dans [22], un décodeur LS-binaural adaptatif a été proposé, qui constraint la matrice de covariance de l'audio décodé à refléter celle de la matrice de covariance des HRTF des auditeurs. La méthode ne nécessite pas d'estimation explicite de l'équilibre direct/diffus, ni n'a besoin d'utiliser la décorrélation.

### *1.2.1 Décodage binaural paramétrique avec mixage optimal*

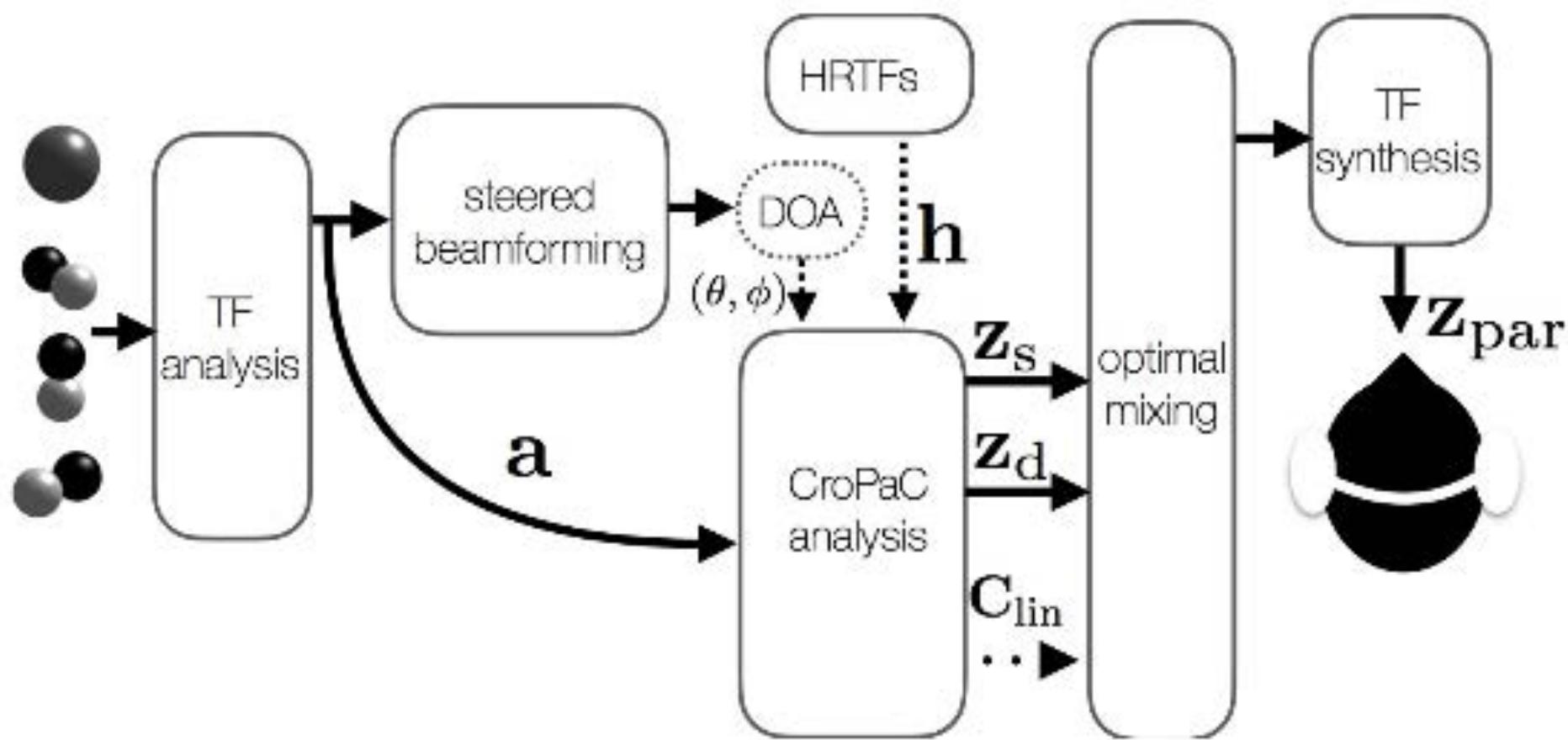
Le principal point de critique concernant les méthodes de reproduction paramétrique est l'incursion d'artefacts variant dans le temps. Celles-ci se produisent soit parce que la scène d'entrée n'est pas conforme au modèle de champ sonore supposé, soit parce que les techniques de rendu ne sont pas suffisamment robustes. Par conséquent, dans [28], une technique de mélange optimale a été proposée, qui tente de synthétiser la sortie en utilisant une combinaison linéaire de signaux prototypes décodés linéairement, autant que possible ; Donc en conservant une grande partie de la haute fidélité monocanal dans la sortie. Le procédé n'emploie alors la décorrélation que si les dépendances inter-canaux de sortie s'écartent toujours de la cible, atténuant ainsi les éventuels artefacts de décorrélation ; comme le maculage temporel des transitoires.

L'approche est formulée dans le domaine de la covariance et repose sur la construction de matrices de covariance cibles à bande étroite variant dans le temps, qui sont dictées par le modèle de champ sonore de la méthode paramétrique. Il utilise ensuite l'audio décoded linéairement traditionnel comme prototype. L'approche a été utilisée précédemment dans [19] en utilisant le modèle de champ sonore DirAC, et aussi dans la solution de forme fermée de [22] pour la reproduction binaurale. Cependant, en principe, n'importe quel modèle de champ sonore peut utiliser cette approche de mixage optimale pour l'étage de synthèse.

## Motivation pour ce travail

La poursuite du développement des décodeurs ambisoniques du premier ordre présente un intérêt particulier ; compte tenu de la prévalence des réseaux de microphones ambisoniques de premier ordre disponibles dans le commerce et du matériel accessible sur Internet. Des progrès récents ont montré des améliorations dans la performance perçue du décodage ambisonique binaural linéaire [13, 15]. Cependant, ces décodeurs reposent toujours sur des ordres d'entrée qui peuvent être considérés comme peu pratiques pour une adoption à grande échelle aujourd'hui. Surtout compte tenu de la croissance quadratique du nombre de canaux (et de microphones) avec un ordre croissant, conduisant à des réseaux plus coûteux et à des bandes passantes plus élevées. Par conséquent, le besoin d'alternatives paramétriques robustes dépendant du signal pour les ordres inférieurs est bien établi.

Dans cet article, un décodeur binaural de premier ordre est proposé 1 , qui utilise un modèle de champ sonore paramétrique qui suppose une source et une composante ambiante directionnelle par pavé temps-fréquence. Le modèle est inspiré du modèle ambiant multi- sources et non isotrope utilisé par la méthode COMPASS [25]. Cependant, étant donné la faible résolution des entrées de premier ordre, la méthode proposée : Synoptique général de la méthode proposée emploie un post-filtre spatial supplémentaire CroPaC [4] pour mieux isoler la composante source, et un post-filtre CroPaC inverse pour obtenir la composante ambiante. La méthode proposée extrait également les estimations instantanées de la direction d'arrivée (DoA) via des cartes d'activité de puissance de réponse dirigée (SRP) et la recherche de pics ; renonçant ainsi au besoin d'une moyenne temporelle longue des matrices de covariance d'entrée, comme l'exigent les alternatives de sous-espace. La méthode utilise également la solution de mélange optimale de [5], afin de minimiser la quantité d'énergie de signal décorrélée dans la sortie, et de synthétiser les dépendances inter- canaux cibles de manière linéaire, autant que possible.



Synoptique général de la méthode proposée.

## ÉVALUATION

Un test multistimulus a été réalisé, afin de comparer les performances perçues de l'approche proposée (CroPaC1), avec les décodeurs ambisoniques à rééchantillonnage spatial du premier ordre (Ambi1) et du troisième ordre (Ambi3) [12]. Notez que l'approche de mélange optimal (12) a été utilisée, en utilisant le même décodeur ambisonique pour le prototype zlin.

Des scènes de test synthétiques ont été créées comme suit : plusieurs haut-parleurs (haut-parleurs), un groupe de funk moderne (groove) et un mix comprenant un haut-parleur, des applaudissements, une fontaine à eau et un piano (mix). Les sources individuelles ont été binaurisées directement pour les cas de test de référence, et codées en signaux de premier et troisième ordre et passées par leurs décodeurs respectifs pour les cas de test secs. Un simulateur de salle de source d'image boîte à chaussures 2 a été utilisé pour obtenir des contreparties réverbérantes (rev). Le simulateur de salle a été configuré pour ressembler à une petite salle. Les sources d'images arrivant à la position du récepteur ont été directement binaurales pour les cas de référence et codées en signaux harmoniques sphériques pour les cas de test réverbérants.

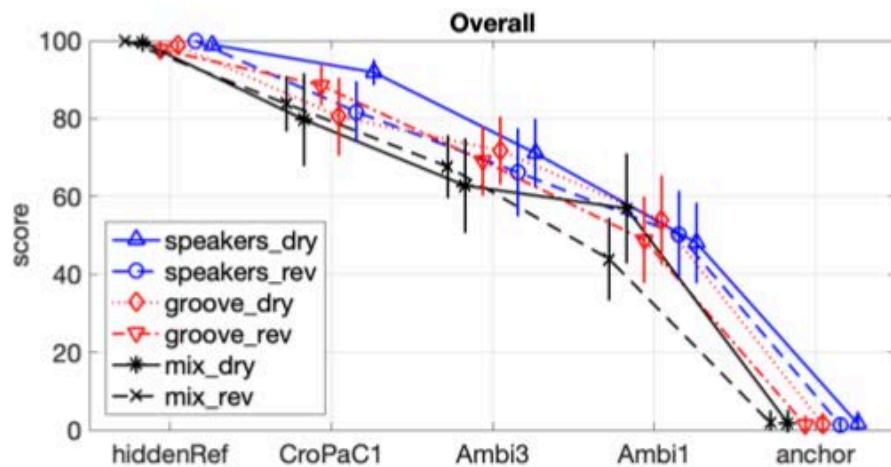
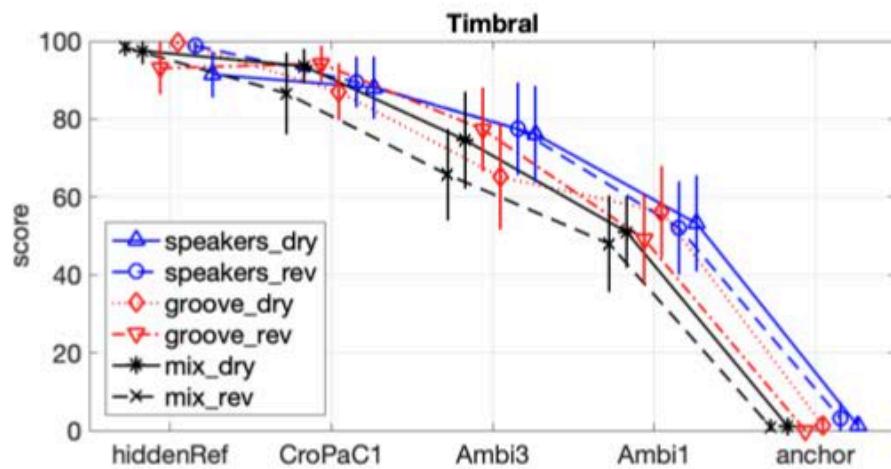
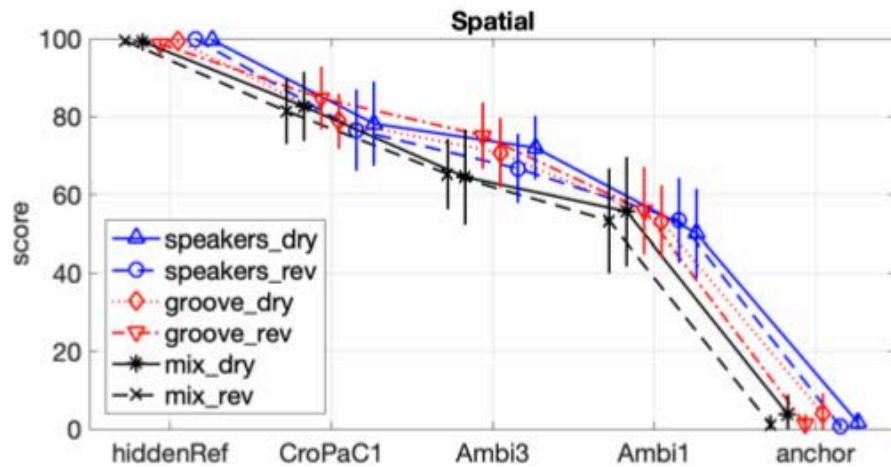
L'évaluation comportait trois parties, qui traitaient des différences spatiales, timbrales et globales entre les méthodes. Pour l'évaluation des différences spatiales, les spectres moyens de la référence ont été imposés aux cas tests. Par

conséquent, les différences de timbre entre les méthodes ont été grandement atténuées, mais ont conservé leurs caractéristiques spatiales d'origine. Les sujets de test ont été explicitement instruits d'ignorer toutes les différences de timbre restantes, qui peuvent ne pas avoir été traitées par l'égalisation. L'ancre a été obtenue sous la forme d'une composante harmonique sphérique omnidirectionnelle, répliquée sur chaque canal binaural et également égalisée. Pour évaluer les différences de timbre, les spectres moyens de chaque cas de test ont été imposés sur un signal omnidirectionnel et envoyés aux canaux gauche et droit. Par conséquent, les différences spatiales entre les méthodes ont été éliminées, ne retenant ainsi que les différences timbrales.

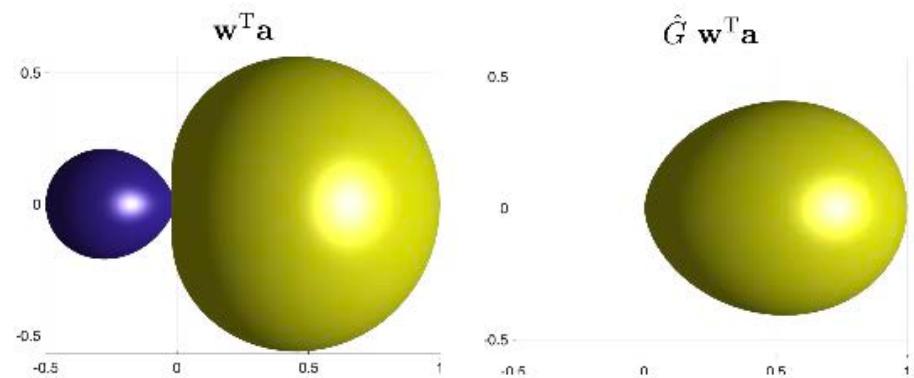
L'ancre a été obtenue comme les spectres moyens de la référence, répliqués sur chaque canal binaural et filtrés passe-bas à 4 kHz. Enfin, pour évaluer les différences globales, seul le RMS large bande de la référence a été utilisé pour normaliser les cas de test. Par conséquent, toutes les différences timbrales et spatiales entre les cas de test sont restées, et les participants au test ont noté les échantillons en fonction de leur pondération subjective des attributs reproduits. Au total, 14 auditeurs experts ont participé aux tests d'écoute dans des cabines d'écoute spécialement conçues. Les auteurs présents n'ont pas participé aux tests (les moyennes et les intervalles de confiance à 95 %). On peut observer que les caractéristiques spatiales et timbrales de la référence ont été reproduites plus fidèlement en utilisant la méthode proposée, par rapport à l'ambisonique conventionnelle avec la même entrée de premier ordre. De plus, la méthode proposée a donné des scores plus conformes à ceux des ambisoniques de troisième ordre. Il convient de souligner que l'ambisonique de troisième ordre utilise quatre fois plus de canaux d'entrée que celui de la méthode proposée. Cela représente donc une réduction significative de la bande passante, sans compromettant la précision ou la fidélité spatiale perçue.

## CONCLUSION

Cet article a proposé un décodeur ambisonique binaural paramétrique de premier ordre, qui utilise un modèle de champ sonore comprenant une source et une composante ambiante directionnelle par pavé temps-fréquence. L'approche proposée isole d'abord les composantes de la source à l'aide d'un formateur de faisceau de domaine harmonique sphérique, modulé par le post-filtre spatial Cross-Pattern Coherence (CroPaC). Le flux ambiant est alors simplement le résidu, une fois que les composantes de la source ont été soustraites du champ sonore d'entrée. L'approche proposée est inspirée de la méthode COMPASS. Cependant, avec le post-filtre CroPaC, il utilise également une estimation instantanée de la direction de la source et synthétise la sortie de manière linéaire autant que possible ; afin d'améliorer la fidélité des signaux de sortie.



Les moyennes et les intervalles de confiance à 95 % pour chaque scène sonore individuelle. Les critères d'évaluation étaient : spatial uniquement, timbral uniquement et global (haut-bas).



Un formateur de faisceau hyper cardioïde de premier ordre sans (à gauche) et avec (à droite) le post-filtre CroPaC.

Merci de votre attention

Site : <https://www.lesonbinaural.fr>

Mail : [b.lagnel@gmail.com](mailto:b.lagnel@gmail.com)