

Optimized Covariance Domain Framework for Time-Frequency Processing of Spatial Audio

JUHA VILKAMO^{1,2}, TOM BÄCKSTRÖM², AND ACHIM KUNTZ²

¹*Aalto University, Espoo, Finland*

²*Fraunhofer IIS, Erlangen, Germany*

The covariance matrix of a multichannel audio signal is a measure that contains the channel energies and the inter-channel dependencies. This measure is perceptually relevant in frequency bands, since with the effect of the acoustic transfer path, it forms the inter-aural cues based on which the human spatial hearing system decodes the spatial sound. This is the foundation to several state of the art perceptual spatial processing techniques. In this paper we propose a generalized framework for spatial sound processing that operates directly in the covariance matrix domain. The technique optimizes the sound quality by a least mean square solution with respect to a defined prototype signal. Furthermore, the technique uses all available independent signal components in the input channels in order to minimize the undesired effects due to the decorrelators that are often employed in the concerned use cases. A wide variety of applications of the proposed method has been identified.

1 INTRODUCTION

Our human auditory system is able to localize sound sources, and in some degree estimate the size and characteristics of the surrounding space, based on the two audio signals entering the ear canals. The acoustic effect of the head, torso, and pinnae causes the sound to attain characteristics that are specific to the angle of the arrival. These are the inter-aural level-difference (ILD) and the inter-aural phase difference (IPD) as function of the frequency [1]. In addition, the inter-aural coherence (IC) becomes relevant in the presence of several sources or in diffuse sound environments such as reverberation.

In perceptual time-frequency processing of spatial audio [2,3,4,5,6,7,8], the processing focuses on those characteristics of the stereo or multichannel sound that contribute to the above binaural cues. These are in frequency bands the inter-channel level-difference (ICLD), inter-channel phase difference (ICPD), and the inter-channel coherence (ICC) [9]. In binaural reproduction over headphones, the inter-channel cues are the same as the inter-aural cues. In loudspeaker reproduction, the inter-channel cues map to the binaural cues with the effect of the acoustic transfer path from the loudspeakers to the listeners' ears.

Several techniques based on the inter-channel cues have emerged. Efficient multichannel transmission systems [2,3]

transmit only a reduced number of channels and a low-bitrate parametric side information with which the original inter-channel cues and thus the original spatial perception are recovered. Spatial Audio Object Coding [4] is a similar technique that transmits several audio objects in down-mix channels and allows an option to render the objects flexibly in the receiver end. Directional Audio Coding [5] is a DSP-driven microphone technique that analyzes the sound field direction and diffuseness in frequency bands, and uses amplitude panning and decorrelation to synthesize such a multi-channel loudspeaker signal that produces a sound field perceptually close to the original. Panning and decorrelation are tools for generating a set of ICLDs and ICCs that contribute in generating the desired spatial perception. Often such applications apply the processing in a non-uniform frequency resolution, such as the equivalent rectangular bandwidth (ERB) [10] scale.

In this paper we propose a generalized and optimized framework for spatial audio processing in terms of the signal covariance matrix in frequency bands, which is a measure that contains the described inter-channel characteristics. Applications that apply the framework determine adaptively the target covariance matrix in the time-frequency domain, and the framework performs the transform while aiming to best preserve the overall sound quality.

In addition to the covariance matrix, the signal content in the channels also matter. For example, in a scenario of teleconferencing with spectrally overlapping talkers, it is

Correspondence should be addressed to juha.vilkamo@aalto.fi

relevant to reproduce the harmonics of each talker approximately at their correct directions. Similarly, in reproduction of sound recorded in echoic spaces, due to the precedence effect [1,11] it is important to reproduce the first arriving wavefront at its correct direction. For this reason the spatial synthesis in the proposed framework is controlled also by a prototype signal. The proposed framework produces an output signal that is as close as possible to the prototype signal, with the boundary condition that the target covariance matrix, i.e., the target spatial characteristic, is attained.

The proposed framework also optimizes the usage of the decorrelators, which are processes that provide new incoherent signal components based on the input signals. By definition, they affect the time- and phase structure of the signal, which may be detrimental to the perceived signal quality. The proposed framework ensures that the decorrelated energy is injected only to the minimum necessary extent, i.e., only to complement for the missing independent signal components in the input channels.

The main expected application of the proposed framework is in the field of spatial sound recording [5,7], in which the problems related to the signal covariance are particularly apparent due to the physical limitations of building high-quality microphones that are directional enough for spatial sound reproduction. Further expected use cases include stereo and multichannel enhancement, ambience extraction, upmixing, and downmixing. The proposed framework can operate with real or complex definitions of the covariance matrices. The real-valued definitions can be sufficient in applications where the phase differences are less relevant, e.g., in spatial sound rendering using a set of coincident microphones. Complex valued processing is necessary in applications that rely also on phase differences, such as binaural rendering and spatial sound rendering from a spaced microphone array.

The paper is organized as follows. We start by the definitions, followed by the derivation of the proposed method. The adaptive mixing solution is derived first using a set of simplifying temporary assumptions. Then the generalization is provided and the means to apply a decorrelated residual signal to compensate when there is not enough independent signal energy available otherwise. Then an example implementation in context of stereo upmixing is given, followed by a set of further numerical examples, and finally the conclusion. As an appendix, an example Matlab implementation providing the optimized adaptive mixing solution as described in this paper is provided.

2 DEFINITIONS

Matrices and vectors are denoted with bold faced symbols, where uppercase denotes a matrix. Matrix and vector elements are marked with row and column subindices i and j , such that g_{ij} is an element of matrix \mathbf{G} , and x_i is an element of vector \mathbf{x} . The mixing solution is first derived assuming the same number of the input and output channels N . The method is generalized to arbitrary channel numbers in Section 3.3.

Let us define the input channels $x_i(t)$, $i = 1, \dots, N$ and the output channels $y_j(t)$, $j = 1, \dots, N$. Their short-time Fourier transformed (STFT) counterparts are $X_i(k, l)$ and $Y_j(k, l)$, where k is the downsampled time index and l is the frequency index. The time-frequency signal vectors are

$$\mathbf{x}(k, l) = \begin{bmatrix} X_1(k, l) \\ X_2(k, l) \\ \vdots \\ X_N(k, l) \end{bmatrix} \quad \mathbf{y}(k, l) = \begin{bmatrix} Y_1(k, l) \\ Y_2(k, l) \\ \vdots \\ Y_N(k, l) \end{bmatrix}. \quad (1)$$

Note that the time and frequency indices (k, l) are now omitted for brevity of notation. The signal covariance matrices are

$$\begin{aligned} \mathbf{C}_x &= \mathbf{E}[\mathbf{x}\mathbf{x}^H] \\ \mathbf{C}_y &= \mathbf{E}[\mathbf{y}\mathbf{y}^H], \end{aligned} \quad (2)$$

where $\mathbf{E}[\]$ is the expectation operator and \mathbf{x}^H and \mathbf{y}^H are the conjugate transposes of \mathbf{x} and \mathbf{y} . The target covariance matrix \mathbf{C}_y depends on the specific application, of which an example is given in Section 4.2. Based on \mathbf{C}_y the proposed algorithm applies a mixing procedure to the input signals \mathbf{x} that generates the optimal output signals \mathbf{y} .

In an actual implementation the expectation operator in Eq. (2) is replaced with a mean operator over a number of time and/or frequency indices in a resolution that approximates that of the human spatial hearing. Such a time-frequency area is denoted as the time-frequency tile. The adaptive mixing solution in this paper is formulated once per each tile.

Since the covariance matrices are Hermitian and positive-semi-definite they can be decomposed as

$$\begin{aligned} \mathbf{C}_x &= \mathbf{K}_x \mathbf{K}_x^H \\ \mathbf{C}_y &= \mathbf{K}_y \mathbf{K}_y^H. \end{aligned} \quad (3)$$

Such decompositions can be obtained, e.g., by using the Cholesky decomposition or eigendecomposition [12]. The decomposition is not unique. For any unitary matrices \mathbf{P}_x and \mathbf{P}_y , matrices $\mathbf{K}_x \mathbf{P}_x$ and $\mathbf{K}_y \mathbf{P}_y$ also fulfill the decomposition condition since

$$\begin{aligned} \mathbf{K}_x \mathbf{P}_x \mathbf{P}_x^H \mathbf{K}_x^H &= \mathbf{K}_x \mathbf{K}_x^H = \mathbf{C}_x \\ \mathbf{K}_y \mathbf{P}_y \mathbf{P}_y^H \mathbf{K}_y^H &= \mathbf{K}_y \mathbf{K}_y^H = \mathbf{C}_y. \end{aligned} \quad (4)$$

Finally, let us define the prototype signal $\hat{\mathbf{y}}$, which is applied in the derivation of the least mean square mixing solution. The prototype signal is constructed by applying a prototype matrix \mathbf{Q} to the input signal

$$\hat{\mathbf{y}} = \mathbf{Q}\mathbf{x}. \quad (5)$$

The matrix \mathbf{Q} can be defined static or time or frequency variant, depending on the application and is discussed further in the end of this section. The covariance matrix of the prototype signal is $\mathbf{C}_{\hat{\mathbf{y}}} = \mathbf{E}[\hat{\mathbf{y}}\hat{\mathbf{y}}^H]$, and typically $\mathbf{C}_{\hat{\mathbf{y}}} \neq \mathbf{C}_y$. The proposed adaptive mixing solution produces an output

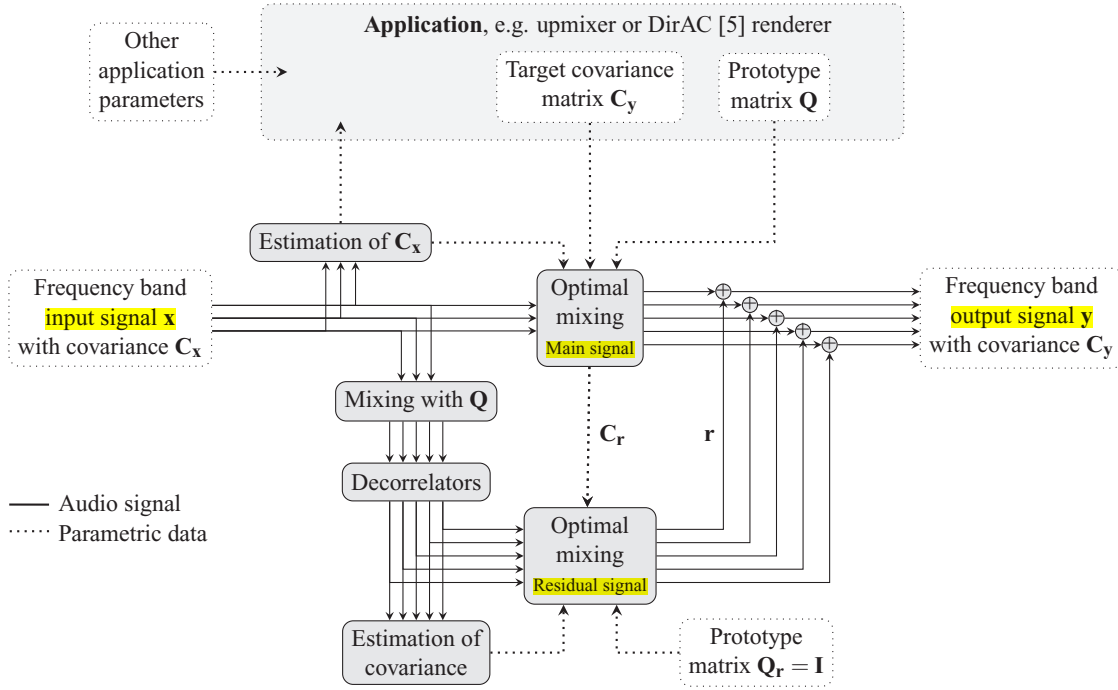


Fig. 1. The proposed framework for perceptual spatial audio processing in terms of the signal covariance matrix in frequency bands. If there are enough independent signal components in the input, the main optimal mixing block is sufficient for generating the output signal with covariance matrix \mathbf{C}_y . The optimal mixing block ensures that the output has the minimum distance to a prototype signal $\hat{\mathbf{y}} = \mathbf{Q}\mathbf{x}$ after its energy normalization. In case there is not enough independent signal energy, a residual signal \mathbf{r} with covariance matrix \mathbf{C}_r that fully compensates for the missing signal components is generated from the decorrelated signals.

signal \mathbf{y} that has \mathbf{C}_y , and minimizes the error

$$e = E[\|\mathbf{y} - \mathbf{G}_{\hat{\mathbf{y}}}\hat{\mathbf{y}}\|^2], \quad (6)$$

where $\mathbf{G}_{\hat{\mathbf{y}}}$ is a diagonal gain matrix that normalizes the per-channel energies of $\hat{\mathbf{y}}$ to those of \mathbf{y} . The normalization step ensures that the error measure is weighted with the channel energies of \mathbf{y} , which is not guaranteed if the normalization is omitted. The diagonal elements of $\mathbf{G}_{\hat{\mathbf{y}}}$ are

$$g_{\hat{\mathbf{y}}_{ii}} = \sqrt{\frac{c_{y_{ii}}}{c_{\hat{\mathbf{y}}_{ii}}}}, i = 1, \dots, N. \quad (7)$$

The general guideline for designing the application specific prototype matrix \mathbf{Q} is to consider which kind of signal content is appropriate for each output channel. In spatial sound recording, for example, it is desirable that the sound emitted by each of the loudspeakers is most similar to the sound that in the recorded sound scene arrived from the corresponding angle. Thus, such \mathbf{Q} is selected that mixes the microphone signals in the way that a spatial amplification pattern is directed toward the desired angle, while other angles are attenuated. As a further example, in stereo upmixing a reasonable \mathbf{Q} is such that defines that the loudspeaker channels at one side correspond to the original stereo channel at the same side.

3 ADAPTIVE MIXING SOLUTION

The objective is to generate from the input signal \mathbf{x} an output signal \mathbf{y} that has the covariance matrix \mathbf{C}_y , while the error measure in Eq. (6) is minimized. Two processing methods are utilized in the following. The preferred method is adaptive mixing of the input channels since it best preserves the signal quality. The complementary method is the injection of a decorrelated residual signal when there are not enough independent signal components available otherwise. The input-output relation can be written as

$$\mathbf{y} = \mathbf{M}\mathbf{x} + \mathbf{r}, \quad (8)$$

where \mathbf{M} is a mixing matrix and \mathbf{r} is the residual signal processed from the decorrelated signals. The overall block diagram of the processing is shown in Fig. 1.

3.1 Adaptive Mixing Solution in Idealized Condition

In addition to the temporary assumption of the same number of the input and output channels, let us also assume that the input signal is full rank. Consequently, as will be shown, the residual signal is not necessary, i.e., $\mathbf{r} = \mathbf{0}$. Let us use a notation $\tilde{\mathbf{M}}$ in place of \mathbf{M} to denote the adaptive mixing solution achieved with the temporary assumptions. From the simplified signal relation $\mathbf{y} = \tilde{\mathbf{M}}\mathbf{x}$ and the definition of the covariance matrices in Eq. (2), it follows that

$$\mathbf{C}_y = E[\mathbf{y}\mathbf{y}^H] = E[\tilde{\mathbf{M}}\mathbf{x}\mathbf{x}^H\tilde{\mathbf{M}}^H] = \tilde{\mathbf{M}}\mathbf{C}_x\tilde{\mathbf{M}}^H. \quad (9)$$

When the covariance matrices in Eq. (9) are decomposed according to Eq. (4) it follows that

$$\mathbf{K}_y \mathbf{P}_y \mathbf{P}_y^H \mathbf{K}_y^H = \tilde{\mathbf{M}} \mathbf{K}_x \mathbf{P}_x \mathbf{P}_x^H \mathbf{K}_x^H \tilde{\mathbf{M}}^H, \quad (10)$$

from which the set of solutions for $\tilde{\mathbf{M}}$ that fulfill Eq. (9) follows

$$\tilde{\mathbf{M}} = \mathbf{K}_y \mathbf{P}_y \mathbf{P}_x^H \mathbf{K}_x^{-1} = \mathbf{K}_y \mathbf{P} \mathbf{K}_x^{-1}. \quad (11)$$

The condition for these solutions is that \mathbf{K}_x^{-1} exists. The unitary matrix $\mathbf{P} = \mathbf{P}_y \mathbf{P}_x^H$ is the remaining free parameter. In Appendix A it is shown that the error measure in Eq. (6) is minimized and thus an optimal adaptive mixing solution is achieved when

$$\mathbf{P} = \mathbf{V} \mathbf{U}^H, \quad (12)$$

where the unitary matrices \mathbf{V} and \mathbf{U} are from the singular value decomposition $\mathbf{U} \mathbf{S} \mathbf{V}^H = \mathbf{K}_x^H \mathbf{Q}^H \mathbf{G}_y^H \mathbf{K}_y$, and \mathbf{S} is a non-negative diagonal matrix.

3.2 Adaptive Mixing Solution Generalized

In practice it is possible, and in some applications common, that \mathbf{K}_x^{-1} in Eq. (11) does not exist, or that this inverse entails very large multipliers if some of the principal components in \mathbf{x} are very small. An effective way to regularize the inverse is by the singular value decomposition $\mathbf{K}_x = \mathbf{U}_x \mathbf{S}_x \mathbf{V}_x^H$. The unregularized inverse is

$$\mathbf{K}_x^{-1} = \mathbf{V}_x \mathbf{S}_x^{-1} \mathbf{U}_x^H. \quad (13)$$

The problems arise when some of the diagonal values of the non-negative diagonal matrix \mathbf{S}_x are zero or very small. It is straightforward to manually tune a regularization function in this domain. In the example implementation, we first define $s_{x_{\max}}$ as the largest diagonal element of \mathbf{S}_x and then form the regularized diagonal matrix \mathbf{S}'_x with entries

$$s'_{x_{ii}} = \max [s_{x_{ii}}, \alpha s_{x_{\max}}], i = 1, \dots, N. \quad (14)$$

where α is a factor between 0 and 1 that determines the amount of regularization. In the informal tests the value $\alpha = \frac{1}{5}$ was found a good compromise that avoids amplifying small energy signal components too much, but still keeping a good ability for the framework to take benefit of the available independent signal energy in the input channels. The regularized inverse is $\mathbf{K}'^{-1}_x = \mathbf{V}_x \mathbf{S}'^{-1}_x \mathbf{U}_x^H$, and the corresponding regularized mixing matrix is $\mathbf{M} = \mathbf{K}_y \mathbf{P} \mathbf{K}'^{-1}_x$. The regularization effectively means that within the mixing process the amplification of some of the small principal components in \mathbf{x} is reduced, and consequently their impact to the output signal \mathbf{y} is also reduced, and the target covariance \mathbf{C}_y is in general not reached.

The motivation for injecting the residual signal, which is discussed in the following, is to achieve \mathbf{C}_y also in the regularized conditions. The residual signal \mathbf{r} is incoherent with respect to \mathbf{x} , and its covariance matrix $\mathbf{C}_r = \mathbf{E}[\mathbf{r} \mathbf{r}^H]$

can be derived by setting the target covariance condition

$$\begin{aligned} \mathbf{C}_y &= \mathbf{E}[(\mathbf{M}\mathbf{x} + \mathbf{r})(\mathbf{M}\mathbf{x} + \mathbf{r})^H] \\ &= \mathbf{E}[(\mathbf{M}\mathbf{x}\mathbf{x}^H \mathbf{M}^H] + \mathbf{E}[\mathbf{r}\mathbf{r}^H] \\ &= \mathbf{M} \mathbf{C}_x \mathbf{M}^H + \mathbf{C}_r \\ &\Rightarrow \boxed{\mathbf{C}_r = \mathbf{C}_y - \mathbf{M} \mathbf{C}_x \mathbf{M}^H}. \end{aligned} \quad (15)$$

Any signal that is independent with respect to \mathbf{x} that is processed to have the covariance \mathbf{C}_r , serves as a residual signal that enables reconstructing the target covariance matrix \mathbf{C}_y in situations where the regularization as in Eq. (14) was applied. As illustrated in Fig. 1, such a residual signal \mathbf{r} can be generated by applying decorrelators to the prototype signal $\hat{\mathbf{y}}$ and processing the result using the proposed adaptive mixing solution to attain \mathbf{C}_r . In the residual mixing process, the matrix \mathbf{Q}_r is defined as an identity matrix, since the prototype signal is the decorrelated signal itself. Note that although the residual signal is generated from the decorrelated signals, the residual covariance matrix \mathbf{C}_r is typically not diagonal.

If the usage of decorrelators is to be completely avoided, e.g., in the higher frequency range where the human hearing is phase insensitive, at least the target channel energies can be achieved by multiplying the rows of \mathbf{M} so that

$$\mathbf{M}' = \mathbf{G} \mathbf{M}, \quad (16)$$

where \mathbf{G} is a diagonal gain matrix with elements

$$g_{ii} = \sqrt{\frac{c_{y_{ii}}}{c_{\tilde{y}_{ii}}}}, i = 1, \dots, N, \quad (17)$$

where $c_{\tilde{y}_{ii}}$ is the i th diagonal element of $\mathbf{C}_{\tilde{y}} = \mathbf{M} \mathbf{C}_x \mathbf{M}^H$. Note that although the matrix \mathbf{M}' is no more optimal in a mathematical sense, it is well-warranted in a perceptual sense.

3.3 Different Number of Input and Output Channels

In many applications the number of input and output channels is different. The regularization process presented in Section 3.2 provides robust mixing solutions also when some of the channels are zero-valued. The solution to the different channel numbers is thus directly available by assuming that the signal with fewer channels is extended to the higher dimension with additional zero-valued channels.

Such an approach however implies computational overhead since some rows or columns of the mixing matrix \mathbf{M} correspond to the zero-valued channels. The same mixing result is obtained more efficiently by introducing a matrix Λ that is an identity matrix appended with zeros to dimension $N_y \times N_x$, where N_y is the number of the output channels, and N_x is the number of the input channels, e.g.,

$$\Lambda_{3 \times 2} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}. \quad (18)$$

When \mathbf{P} is redefined as

$$\mathbf{P} = \mathbf{V} \Lambda \mathbf{U}^H \quad (19)$$

the resulting \mathbf{M} is a $N_y \times N_x$ mixing matrix that is the same as the relevant part of the \mathbf{M} of the zero-padding case. Consequently, \mathbf{C}_x , \mathbf{C}_y , \mathbf{K}_x , and \mathbf{K}_y can be of their natural dimension, and the prototype matrix \mathbf{Q} is of dimension $N_y \times N_x$.

3.4 Decomposability of the Covariance Matrices

The input signal covariance matrix is always decomposable to $\mathbf{C}_x = \mathbf{K}_x \mathbf{K}_x^H$ because it is a positive semi-definite measure from an actual signal. It is, however, possible to define such target covariance matrices that are not decomposable for the reason that they represent impossible channel dependencies. There are methods to ensure decomposability, such as adjusting the negative eigenvalues to zeros and normalizing the energy [13]. However, the most meaningful usage of the proposed method is to request only valid covariance matrices.

4 EXAMPLE IMPLEMENTATION

In this section it is shown by an example how the proposed framework can be implemented in a practical spatial sound processing application. We design a processing rule for automatic stereo to surround upmixing, in which the direct and ambience components of the stereo signal are estimated and then distributed to the larger set of loudspeakers.

4.1 Prototype Matrix \mathbf{Q} in Upmixing

Assuming 5.0 channel order L, R, C, Ls, Rs, a reasonable prototype matrix in stereo upmixing is

$$\mathbf{Q} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ \sqrt{0.5} & \sqrt{0.5} \\ 1 & 0 \\ 0 & 1 \end{bmatrix}. \quad (20)$$

This definition means that the prototype signal for the two left channels is the original left channel, and similarly for the right, and that the prototype signal for the center channel is the mixture.

4.2 Target Covariance Matrix in Upmixing

Following the principles in [6], the direct-ambience model parameters describing each stereo time-frequency tile can be obtained: the energy of the panned direct component E_D , the energy of the incoherent ambience component E_A , and the panning angle of the direct component α_D . The target covariance matrix for each time-frequency tile is built as function of these parameters. A reasonable rule is to distribute the ambience energy to incoherently surround the listener and to amplitude pan the direct component using a pair from the three frontal channels.

A practical choice for distributing the ambience energy is incoherently using the channels L, R, Ls, and Rs. The center channel is omitted since it is assumed not to contribute much to the perceived spatial quality of the ambience elements. This procedure may improve the overall

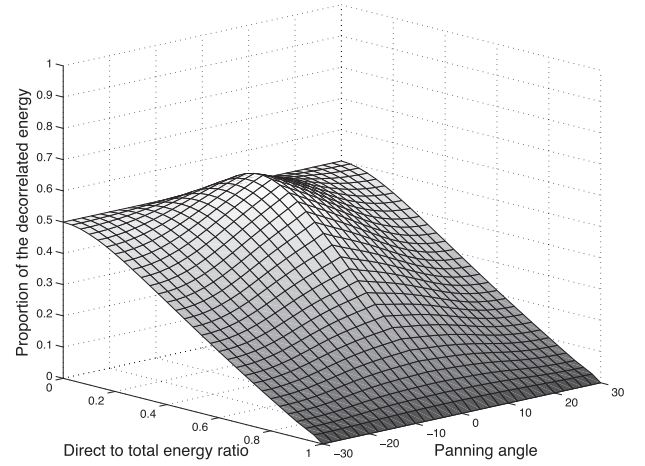


Fig. 2. The proportion of the output energy that is from the decorrelators with the example upmixing implementation.

sound quality with the proposed method since having fewer independent output channels entails lesser relative amount of the decorrelated signal energy. According to this design principle, a covariance matrix describing the diffuse sound energy is

$$\mathbf{C}_A = E_A \begin{bmatrix} 0.25 & 0 & 0 & 0 & 0 \\ 0 & 0.25 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0.25 & 0 \\ 0 & 0 & 0 & 0 & 0.25 \end{bmatrix}. \quad (21)$$

For the direct part, let us assume a column vector \mathbf{v} that contains the vector base amplitude panning (VBAP) [14] gains corresponding to the panning angle α_D . Vector \mathbf{v} has only two non-zero values for the pair of the front loudspeakers that is active in the panning. The direct part covariance matrix is

$$\mathbf{C}_D = E_D \mathbf{v} \mathbf{v}^T. \quad (22)$$

Finally, since the direct and ambient parts were defined incoherent with respect to each other, the overall target covariance matrix is

$$\mathbf{C}_y = \mathbf{C}_A + \mathbf{C}_D. \quad (23)$$

4.3 Simulations—Applied Amount of the Decorrelated Sound Energy

A set of test scenarios was generated in Matlab by varying the panning angle and the direct-to-total energy ratio $E_D/(E_D + E_A)$ of a simulated stereo mix. In all conditions, the target covariance matrix was formulated using the example implementation in Appendix B, with the control matrices as in Eqs. (20)–(23). Fig. 2 shows the resulting relative amount of the decorrelated energy in the output signal in all simulated conditions. The maximum value 3/5 is reached at 0° semi-ambient condition, since in this condition the applied rule for the target covariance matrix generates a diagonal equal energy covariance matrix, which is the extreme in terms of the requested independent signal energy. In other conditions the target covariance matrix becomes

Table 1. Numerical examples of the proposed algorithm. When a signal with covariance is processed with a mixing matrix, and completed with a possible residual signal with, the output signal has the covariance. Although these numerical examples are static, the typical use case of the proposed method is dynamic. The channel order is assumed $L, R, C, Ls, Rs, (Lr, Rr)$.

Context	Input parameter matrices			Formulated processing matrices	
	C_x	Q	C_y	M	C_r
Decorrelation: High input ICC	$\begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$	$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$	$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$	$\begin{bmatrix} 1.5 & -0.75 \\ -0.75 & 1.5 \end{bmatrix}$	$\begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$
Decorrelation: Very high input ICC	$\begin{bmatrix} 1 & 0.97 \\ 0.97 & 1 \end{bmatrix}$	$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$	$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$	$\begin{bmatrix} 2.1 & -1.4 \\ -1.4 & 2.1 \end{bmatrix}$	$\begin{bmatrix} 0.31 & -0.31 \\ -0.31 & 0.31 \end{bmatrix}$
Stereo upmixing	$\begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$	$\begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0.71 & 0.71 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}$	$\begin{bmatrix} 0.5 & 0 & 0 & 0 & 0 \\ 0 & 0.5 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 \\ 0 & 0 & 0 & 0.5 & 0 \\ 0 & 0 & 0 & 0 & 0.5 \end{bmatrix}$	$\begin{bmatrix} 0.31 & -0.19 \\ -0.19 & 0.31 \\ 0.53 & 0.53 \\ 0.31 & -0.19 \\ -0.19 & 0.31 \end{bmatrix}$	$\begin{bmatrix} 0.36 & 0.11 & -0.18 & -0.14 & 0.11 \\ 0.11 & 0.36 & -0.18 & 0.11 & -0.14 \\ -0.18 & -0.18 & 0.32 & -0.18 & -0.18 \\ -0.14 & 0.11 & -0.18 & 0.36 & 0.11 \\ 0.11 & -0.14 & -0.18 & 0.11 & 0.36 \end{bmatrix}$
5.0 to 7.0 upmixing	$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}$	$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0.71 & 0 \\ 0 & 0 & 0 & 0.71 & 0 \\ 0 & 0 & 0 & 0 & 0.71 \end{bmatrix}$	$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0.5 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0.5 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0.5 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0.5 \end{bmatrix}$	$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0.5 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0.5 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0.5 \end{bmatrix}$	$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0.25 & 0 & -0.25 & 0 \\ 0 & 0 & 0 & 0 & 0.25 & 0 & -0.25 \\ 0 & 0 & 0 & -0.25 & 0 & 0.25 & 0 \\ 0 & 0 & 0 & 0 & -0.25 & 0 & 0.25 \end{bmatrix}$
Downmixing: with some non-zero coherences	$\begin{bmatrix} 1 & 0 & 0 & 0.5 & 0 \\ 0 & 1 & 0 & 0 & 0.5 \\ 0 & 0 & 1 & 0 & 0 \\ 0.5 & 0 & 0 & 1 & 0 \\ 0 & 0.5 & 0 & 0 & 1 \end{bmatrix}$	$\begin{bmatrix} 1 & 0 & 0.71 & 1 & 0 \\ 0 & 1 & 0.71 & 0 & 1 \end{bmatrix}$	$\begin{bmatrix} 2.5 & 0.5 \\ 0.5 & 2.5 \end{bmatrix}$	$\begin{bmatrix} 0.84 & 0.02 & 0.61 & 0.84 & 0.02 \\ 0.02 & 0.84 & 0.61 & 0.02 & 0.84 \end{bmatrix}$	$\begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$
Decorrelation: Standard 5.0 layout coincident hypercardioids in a diffuse field	$\begin{bmatrix} 1 & 0.65 & 0.91 & 0.43 & -0.22 \\ 0.65 & 1 & 0.91 & -0.22 & 0.43 \\ 0.91 & 0.91 & 1 & 0.07 & 0.07 \\ 0.43 & -0.22 & 0.07 & 1 & -0.22 \\ -0.22 & 0.43 & 0.07 & -0.22 & 1 \end{bmatrix}$	$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}$	$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}$	$\begin{bmatrix} 2 & -0.51 & -0.83 & -0.53 & 0.41 \\ -0.51 & 2 & -0.83 & 0.41 & -0.53 \\ -0.83 & -0.83 & 2.1 & 0.04 & 0.04 \\ -0.53 & 0.41 & 0.04 & 1.2 & -0.07 \\ 0.41 & -0.53 & 0.04 & -0.07 & 1.2 \end{bmatrix}$	$\begin{bmatrix} 0.58 & -0.2 & -0.34 & -0.23 & 0.19 \\ -0.2 & 0.58 & -0.34 & 0.19 & -0.23 \\ -0.34 & -0.34 & 0.62 & 0.03 & 0.03 \\ -0.23 & 0.19 & 0.03 & 0.11 & -0.11 \\ 0.19 & -0.23 & 0.03 & -0.11 & 0.11 \end{bmatrix}$

sparser and thus there is less necessity for injecting decorrelated energy. In all conditions, the covariance rendering method ensures that the output signal is the closest to the prototype signal.

4.4 Informal Listening

The described upmixer was implemented for automatic processing of stereo music and movie recordings. The applied decorrelators were the band-wise pseudo-random delays as described in [15]. In informal listening the effect was the expected: The center amplitude panned sources were reproduced over the center loudspeaker, the ambience sounds such as the reverberation were expanded to surround the listener, and the overall sound quality was similar to that of the original stereo mix.

5 FURTHER NUMERICAL EXAMPLES

Table 1 shows a set of numerical examples to illustrate the behavior of the proposed method in different use cases. Note that for simplicity, only examples with real-valued covariance matrices were selected. **The mixing matrices M and the residual matrices C_r were formulated with the Matlab code provided in the appendix based on the parameter matrices C_x , C_y , and Q .** Although the matrices are illustrated static, **in typical applications they vary in time and frequency.**

The first and the second row of the table illustrate a use case of **stereo enhancement by means of adjusting the inter-channel coherence to zero**. In the first row there is a small but reasonable incoherent component between the two channels, and thus a fully incoherent output is achieved with only channel mixing. In the second row, the input coherence is very high, i.e., the smaller principal component

is very small. Amplifying this component in extreme degree is not desirable, and thus the built-in limiter starts to require the injection of the decorrelated energy instead, i.e., C_r is now non-zero.

The third row shows a case of stereo to 5.0 upmixing as described in Section 4. **The residual signal is again non-zero since the dimension of the signal is increased.**

The fourth row shows a case of simple 5.0 to 7.0 upmixing, where the original two rear channels are upmixed to the four new rear channels, incoherently. This example illustrates that the processing focuses on those channels where adjustments are requested and leaves the other channels unaffected.

The fifth row depicts a case of downmixing a 5.0 signal to stereo. Passive downmixing, such as applying a static downmixing matrix Q , would amplify the coherent components over the incoherent components. Here the target covariance matrix was defined to preserve the energy, which is fulfilled by the resulting M .

The sixth row illustrates the use case of coincident microphone spatial sound recording. The input covariance matrix C_x is the result of placing ideal hypercardioid microphones to an ideal diffuse field, facing toward the standard angles of the 5.0 setup. The large off-diagonal values in C_x illustrate an inherent disadvantage of passive first order coincident microphone techniques. In the ideal case, the covariance matrix best representing a diffuse field is diagonal, which was therefore set as the target. The relative amount of the resulting decorrelated energy in the output signal is exactly 2/5. This is because there are three independent signal components available in the first order horizontal coincident microphone signals, and two are to be added in order to reach the five-channel diagonal target covariance matrix.

6 CONCLUSION

In this paper we presented a generalized and optimized framework for time-frequency processing of spatial audio in the covariance matrix domain. The method takes into use the available independent signal components in the input channels and thus ensures that the decorrelation is applied only to the minimum necessary extent. The adaptive mixing solution was formulated to maximize the similarity of the output signal and the prototype signal describing the preferred signal content for each of the output channels.

The framework was applied to perform stereo upmixing, and simulations were provided to illustrate the optimized usage of the decorrelated energy in the different input signal conditions. A wide variety of further applications for the proposed framework have been identified.

7 ACKNOWLEDGMENTS

The authors would like to thank Fabian Küch, Giovanni Del Galdo, Emanuël Habets, and Ville Pulkki for their valuable feedback and support.

8 REFERENCES

- [1] J. Blauert, *Spatial Hearing: The Psychophysics of Human Sound Localization* (rev. ed.) (MIT Press, 1997).
- [2] J. Herre, K. Kjörling, J. Breebaart, C. Faller, S. Disch, H. Purnhagen, J. Koppens, J. Hilpert, J. Rödén, W. Oomen, K. Linzmeier, and K. S. Chong, “MPEG Surround—The ISO/MPEG Standard for Efficient and Compatible Multichannel Audio Coding,” *J. Audio Eng. Soc.*, vol. 56, pp. 932–955 (2008 Nov.).
- [3] J. Breebaart, S. van de Par, A. Kohlrausch, and E. Schuijers, “Parametric Coding of Stereo Audio,” *EURASIP J. Applied Signal Processing*, vol. 2005, pp. 1305–1322 (2005).
- [4] J. Herre, H. Purnhagen, J. Koppens, O. Hellmuth, J. Engdegård, J. Hilpert, L. Villemoes, L. Terentiv, C. Falch, A. Hölzer, M. L. Valero, B. Resch, H. Mundt, and H.-O. Oh, “MPEG Spatial Audio Object Coding—The ISO/MPEG Standard for Efficient Coding of Interactive Audio Scenes,” *J. Audio Eng. Soc.*, vol. 60, pp. 655–673 (2012 Sept.).
- [5] V. Pulkki, “Spatial Sound Reproduction with Directional Audio Coding,” *J. Audio Eng. Soc.*, vol. 55, pp. 503–516 (2007 June).
- [6] C. Faller, “Multiple-Loudspeaker Playback of Stereo Signals,” *J. Audio Eng. Soc.*, vol. 54, pp. 1051–1064 (2006 Nov.).
- [7] C. Tournery, C. Faller, F. Küch, and J. Herre, “Converting Stereo Microphone Signals Directly to MPEG-Surround,” presented at the *128th Convention of the Audio Engineering Society* (2010 May), convention paper 7982.
- [8] J. Merimaa and V. Pulkki, “Spatial Impulse Response Rendering I: Analysis and Synthesis,” *J. Audio Eng. Soc.*, vol. 53, pp. 1115–1127 (2005 Dec.).

- [9] F. Baumgarte and C. Faller, “Binaural Cue Coding—Part I: Psychoacoustic Fundamentals and Design Principles,” *Speech and Audio Processing, IEEE Transactions on*, vol. 11, no. 6, pp. 509–519 (2003).
- [10] B. R. Glasberg and B. C. Moore, “Derivation of Auditory Filter Shapes from Notched-Noise Data,” *Hearing Research*, vol. 47, no. 1, pp. 103–138 (1990).
- [11] R. Y. Litovsky, H. S. Colburn, W. A. Yost, and S. J. Guzman, “The Precedence Effect,” *J. Acous. Soc. Am.*, vol. 106, p. 1633 (1999).
- [12] G. H. Golub and C. F. Van Loan, *Matrix Computations*, vol. 3 (JHUP, 2012).
- [13] R. Rebonato and P. Jäckel, “The Most General Methodology for Creating a Valid Correlation Matrix for Risk Management and Option Pricing Purposes,” *Journal of Risk*, vol. 2, no. 2, pp. 17–28 (1999).
- [14] V. Pulkki, “Virtual Sound Source Positioning Using Vector Base Amplitude Panning,” *J. Audio Eng. Soc.*, vol. 45, pp. 456–466 (1997 June).
- [15] J. Vilkamo, “Spatial Sound Reproduction with Frequency Band Processing of B-Format Audio Signals,” Master’s thesis, Helsinki University of Technology (2008).

APPENDIX A. MINIMIZING $\|\mathbf{G}_y \hat{\mathbf{y}} - \mathbf{y}\|^2$ WITH RESPECT TO UNITARY P

From the set of mixing solutions in Eq. (11), we now derive the optimal solution that minimizes the error measure in Eq. (6). Let us consider a signal $\mathbf{w} = \mathbf{K}_x^{-1} \mathbf{x}$, which has the identity covariance matrix

$$\begin{aligned} E[\mathbf{w}\mathbf{w}^H] &= E[\mathbf{K}_x^{-1} \mathbf{x}\mathbf{x}^H (\mathbf{K}_x^{-1})^H] \\ &= \mathbf{K}_x^{-1} E[\mathbf{x}\mathbf{x}^H] (\mathbf{K}_x^{-1})^H \\ &= \mathbf{K}_x^{-1} \mathbf{C}_x (\mathbf{K}_x^{-1})^H \\ &= \mathbf{I}. \end{aligned} \quad (24)$$

Replacing \mathbf{w} in the input-output relation $\mathbf{y} = \tilde{\mathbf{M}}\mathbf{x}$ yields

$$\mathbf{y} = \tilde{\mathbf{M}}\mathbf{x} = \tilde{\mathbf{M}}\mathbf{K}_x \mathbf{w} = \mathbf{K}_y \mathbf{P} \mathbf{w}. \quad (25)$$

The error measure in Eq. (6) can then be written

$$\begin{aligned} e &= E[\|\mathbf{G}_y \hat{\mathbf{y}} - \mathbf{y}\|^2] \\ &= E[\|\mathbf{G}_y \mathbf{Q} \mathbf{x} - \tilde{\mathbf{M}} \mathbf{x}\|^2] \\ &= E[\|\mathbf{G}_y \mathbf{Q} \mathbf{K}_x \mathbf{w} - \mathbf{K}_y \mathbf{P} \mathbf{w}\|^2] \\ &= E[\|(\mathbf{G}_y \mathbf{Q} \mathbf{K}_x - \mathbf{K}_y \mathbf{P}) \mathbf{w}\|^2] \\ &= E[\mathbf{w}^H (\mathbf{G}_y \mathbf{Q} \mathbf{K}_x - \mathbf{K}_y \mathbf{P})^H (\mathbf{G}_y \mathbf{Q} \mathbf{K}_x - \mathbf{K}_y \mathbf{P}) \mathbf{w}]. \end{aligned} \quad (26)$$

From $E[\mathbf{w}\mathbf{w}^H] = \mathbf{I}$, we can readily show for a matrix \mathbf{A} that $E[\mathbf{w}^H \mathbf{A} \mathbf{w}] = \text{tr}(\mathbf{A})$, which is the matrix trace. It follows that Eq. (26) takes the form

$$e = \text{tr}[(\mathbf{G}_y \mathbf{Q} \mathbf{K}_x - \mathbf{K}_y \mathbf{P})^H (\mathbf{G}_y \mathbf{Q} \mathbf{K}_x - \mathbf{K}_y \mathbf{P})]. \quad (27)$$

For matrix traces, we can readily confirm that

$$\begin{aligned}\text{tr}(\mathbf{A} + \mathbf{B}) &= \text{tr}(\mathbf{A}) + \text{tr}(\mathbf{B}) \\ \text{tr}(\mathbf{A}) &= \text{tr}(\mathbf{A}^H) \\ \text{tr}(\mathbf{P}^H \mathbf{A} \mathbf{P}) &= \text{tr}(\mathbf{A}).\end{aligned}\quad (28)$$

Using these properties, Eq. (27) takes the form

$$\begin{aligned}e &= \text{tr}(\mathbf{K}_x^H \mathbf{Q}^H \mathbf{G}_{\hat{\mathbf{y}}}^H \mathbf{G}_{\hat{\mathbf{y}}} \mathbf{Q} \mathbf{K}_x) + \text{tr}(\mathbf{K}_y^H \mathbf{K}_y) \\ &\quad - 2\text{tr}(\mathbf{K}_x^H \mathbf{Q}^H \mathbf{G}_{\hat{\mathbf{y}}}^H \mathbf{K}_y \mathbf{P}).\end{aligned}\quad (29)$$

Only the last term depends on \mathbf{P} . The optimization problem is thus

$$\mathbf{P} = \arg \min_{\mathbf{P}} e = \arg \max_{\mathbf{P}} [\text{tr}(\mathbf{K}_x^H \mathbf{Q}^H \mathbf{G}_{\hat{\mathbf{y}}}^H \mathbf{K}_y \mathbf{P})]. \quad (30)$$

We can readily show for a non-negative diagonal matrix \mathbf{S} and any unitary matrix \mathbf{P}_s that

$$\text{tr}(\mathbf{S}) \geq \text{tr}(\mathbf{S} \mathbf{P}_s). \quad (31)$$

Thereby, by defining the singular value decomposition $\mathbf{U} \mathbf{S} \mathbf{V}^H = \mathbf{K}_x^H \mathbf{Q}^H \mathbf{G}_{\hat{\mathbf{y}}}^H \mathbf{K}_y$, where \mathbf{S} is non-negative and diagonal and \mathbf{U} and \mathbf{V} are unitary, it follows that

$$\begin{aligned}\text{tr}(\mathbf{S}) &\geq \text{tr}(\mathbf{S} \mathbf{V}^H \mathbf{P} \mathbf{U}) = \text{tr}(\mathbf{U} \mathbf{S} \mathbf{V}^H \mathbf{P} \mathbf{U} \mathbf{U}^H) \\ &= \text{tr}(\mathbf{K}_x^H \mathbf{Q}^H \mathbf{G}_{\hat{\mathbf{y}}}^H \mathbf{K}_y \mathbf{P}),\end{aligned}\quad (32)$$

for any unitary \mathbf{P} . The equality holds for

$$\mathbf{P} = \mathbf{V} \mathbf{U}^H, \quad (33)$$

whereby this choice for \mathbf{P} yields the maximum of $\text{tr}(\mathbf{K}_x^H \mathbf{Q}^H \mathbf{G}_{\hat{\mathbf{y}}}^H \mathbf{K}_y \mathbf{P})$ and the minimum of the error measure in Eq. (6).

APPENDIX B. EXAMPLE MATLAB CODE

The following Matlab code was used in the numerical examples and provides the general functionality of the proposed method.

Listing 1: Matlab implementation of the proposed method.

```

1 function [M,Cr]=formulate_M_and_Cr(Cx,Cy,Q,flag)
2 % flag = 0: Expect usage of residuals
3 % flag = 1: Fix energies instead
4 lambda=eye(length(Cy),length(Cx));
5
6 % Decomposition of Cy
7 [U_Cy,S_Cy]=svd(Cy);
8 Ky=U_Cy*sqrt(S_Cy);
9
10 % Decomposition of Cx
11 [U_Cx,S_Cx]=svd(Cx);
12 Kx=U_Cx*sqrt(S_Cx);
13
14 %SVD of Kx
15 Ux=U_Cx;
16 Sx=sqrt(S_Cx);
17 % Vx = identity matrix
18
19 % Regularization Sx
20 Sx_diag=diag(Sx);
21 limit=max(Sx_diag)*0.2+1e-20;
22 Sx_reg_diag=max(Sx_diag,limit);
23
24 % Formulate regularized Kx^-1
25 Kx_reg_inverse=diag(1./Sx_reg_diag)*Ux';
26
27 % Formulate normalization matrix G_hat
28 Cy_hat_diag=diag(Q*Cx*Q');
29 limit=max(Cy_hat_diag)*0.001+1e-20;
30 Cy_hat_diag=max(Cy_hat_diag,limit);
31 G_hat=diag(sqrt(diag(Cy)./Cy_hat_diag));
32
33 % Formulate optimal P
34 [U,S,V]=svd(Kx'*Q'*G_hat'*Ky);
35 P=V*lambda*U';
36
37 % Formulate M
38 M=Ky*P*Kx_reg_inverse;
39
40 % Formulate residual covariance matrix
41 Cy_tilde = M*Cx*M';
42 Cr=Cy-Cy_tilde;
43
44 % Use energy compensation instead of residuals
45 if flag==1
46     adjustment=diag(Cy)./diag(Cy_tilde + 1e-20);
47     G=diag(sqrt(adjustment));
48     M=G*M;
49     Cr='unnecessary';
50 end

```


THE AUTHORS



Juha Vilkkamo



Tom Bäckström



Achim Kuntz

Juha Vilkkamo studied electrical engineering at former Helsinki University of Technology, now Aalto University, Finland. He received his M.Sc. degree in acoustics and audio signal processing in 2008. He worked between 2008 and 2011 as a researcher for Fraunhofer IIS, Germany, in the fields of binaural technologies and spatial sound reproduction. Currently he is pursuing his D.Sc. in Aalto University, in a research collaboration project with Fraunhofer IIS. His enthusiasm at work is in developing optimized and perceptually motivated audio signal processing solutions. His delight in life is his lovely wife and two daughters.

Tom Bäckström is professor of Speech Coding at the International Audio Laboratories Erlangen, which is a joint research unit between the Friedrich-Alexander University and Fraunhofer IIS, Erlangen, Germany. He received his M.Sc. and D.Sc. (tech.) degrees from Helsinki University of Technology, Finland, in 2001 and 2004, respectively,

and continued working there until 2008. Since then he has worked at the AudioLabs Erlangen, first as a researcher and from 2013 as professor. His research interests include speech and audio coding and perception, speech production, time-frequency analysis, matrix and polynomial algebra.

Achim Kuntz received his diploma degree in electrical engineering from the University of Erlangen-Nuremberg in 2002. He then joined the Telecommunications Laboratory at the same university carrying out studies on spatial sound reproduction, multidimensional signals and wave field analysis techniques, resulting in his doctoral degree in 2008. He is currently a researcher at the Fraunhofer Institute for Integrated Circuits IIS in Erlangen and member of the International Audio Laboratories Erlangen. His current primary field of interest are perceptually motivated signal processing algorithms for audio applications including the acquisition, manipulation, coding, and reproduction of spatial sound.