


# On the Single-Sideband Transform for MVDR Beamformers

Vitor Probst Curtarelli<sup>1,\*</sup> , Israel Cohen<sup>1</sup> 

<sup>1</sup> Andrew and Erna Viterbi Faculty of Electrical and Computer Engineering, Technion–Israel Institute of Technology, Technion City, Haifa 3200003, Israel

\* Correspondence: vitor.c@campus.technion.ac.il

**Abstract:** In order to explore different beamforming applications, this paper investigates the application of the Single-Sideband Transform (SSBT) for constructing a Minimum-Variance Distortionless-Response (MVDR) beamformer in the context of the convolutive transfer function (CTF) model for short-window time-frequency transforms by making use of filter-banks and their properties. Our study aims to optimize the appropriate utilization of SSBT in this endeavor, by examining its characteristics and traits. We address a reverberant scenario with multiple noise sources, aiming to minimize both undesired interference and reverberation in the output. Through simulations reflecting real-life scenarios, we show that employing the SSBT correctly leads to a beamformer that outperforms the one obtained when via the Short-Time Fourier Transform (STFT), while exploiting the SSBT's property of it being real-valued. Two approaches were developed with the SSBT, one naive and one refined, with the later being able to ensure the desired distortionless behavior, which is not achieved by the former.

**Keywords:** Single-sideband transform; MVDR beamformer; Filter-banks; Array signal processing; Signal enhancement.

## 1. Introduction

Beamformers are an important tool for signal enhancement, being employed in a plethora of applications from hearing aids [1] to source localization [2] to imaging [3,4]. Among the possible ways to use such devices is to implement them the time-frequency domain [5], which allows the exploitation of frequency-related information while also dynamically adapting to signal changes over time. The most widely used instruments for time-frequency analysis are transforms, from which the Short-Time Fourier Transform (STFT) [6,7] stands out in terms of spread and commonness. However, alternative transforms can also be employed implemented [8–10], each offering unique perspective and information regarding the signal, possibly leading to different outputs.

Among these alternatives, the Single-Sideband Transform (SSBT) [11,12] is of great interest, given its real-valued frequency spectrum. It has been shown that the SSBT works particularly well with short analysis windows [11]. Therefore, if we use the convolutive transfer function (CTF) model [13] to study the desired signal model, the SSBT can lend itself to be useful, if we think about the beamforming process through the lenses of filter-banks [14,15]. Thus, by applying this transform within this context it is possible to pull off

**Citation:** Curtarelli, V. P.; Cohen, I. On the Single-Sideband Transform for MVDR Beamformers. *Algorithms* **2023**, *1*, 0. <https://doi.org/>

Received:  
Revised:  
Accepted:  
Published:

**Copyright:** © 2024 by the authors. Submitted to *Algorithms* for possible open access publication under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

superior performances than only with the STFT. However, it is important to be aware of the limitations of the transform, in order to properly utilize it to try and achieve better outputs.

Two of the most important goals in beamforming are the minimization of noise in the output signal, and the distortionless-ness of the desired signal, both being achieved by the Minimum-Variance Distortionless-Response (MVDR) beamformer [16,17]. As the MVDR beamformer can be used on the time-frequency domain without restrictions on the transform chosen, it is possible to explore and compare the performance of this filter, when designing it through different time-frequency transforms.

Motivated by this, our paper explores the SSB transform and its application on the subject of beamforming within the context of the CTF model. We propose an approach for the CTF that allows the separation of desired and undesired speech components for reverberant environments, and employ this approach for designing the MVDR beamformer. We also explore the traits and limitations of the SSBT, and how to properly adapt the MVDR beamformer to this new transform's constraints. We show that a beamformer designed using the SSBT can surpass the STFT one, while also conforming to the distortionless constraint.

We organized the paper as follows: in Section 2 we introduce the proposed time-frequency transforms, how they're related and what are their relevant properties; Section 3 the considered signal model in the time domain is presented, and how it is transferred into the time-frequency domain; and in Section 4 we develop a true-MVDR beamformer with the SSBT, taking into account its features. In Section 5 we present and discuss the results, comparing the studied methods and beamformers obtained. ?? concludes this paper.

## 2. STFT and the Single-Sideband Transform

When studying signals and systems, often frequency and time-frequency transforms are used in order to change the signal domain [18], allowing the exploitation of different patterns and informations inherent to the signal.

Given a time-domain signal  $x[n]$ , its Short-time Fourier Transform (STFT) [6,7] is

$$X_{\mathcal{F}}[l, k] = \sum_{n=0}^{K-1} w[n] x[n + l \cdot O] e^{-j2\pi k \frac{(n+l \cdot O)}{K}} \quad (1)$$

where  $w[n]$  is an analysis window of length  $K$ ; and  $O$  is the overlap between windows of the transform, usually  $O = \lfloor K/2 \rfloor$ . Even though the STFT is the most traditionally used time-frequency transform, it isn't the only one available. Thus, exploring different possibilities for such an operation can be useful and lead to interesting results.

The Single-Sideband Transform (SSBT) [11] is one such alternative, being cleverly constructed such that its frequency spectrum is real-valued, without loss of information. The SSB transform of  $x[n]$  is defined as

$$X_{\mathcal{S}}[l, k] = \sqrt{2}\Re \left\{ \sum_{n=0}^{K-1} w[n] x[n + l \cdot O] e^{-j2\pi k \frac{(n+l \cdot O)}{K} + j\frac{3\pi}{4}} \right\} \quad (2)$$

Assuming that  $x[n]$  is real-valued, one advantage of using the STFT is that we only need to work with  $\lfloor (K+1)/2 \rfloor + 1$  frequency bins, given its complex-conjugate behavior. Meanwhile, the SSBT requires all  $K$  bins to correctly capture all information of  $x[n]$ , however it is real-valued.

Assuming that all  $K$  bins of the STFT are available, from Eqs. (1) and (2) we have

$$\begin{aligned} X_S[l, k] &= \sqrt{2} \Re \left\{ X_{\mathcal{F}}[l, k] e^{j \frac{3\pi}{4}} \right\} \\ &= -\Re \{ X_{\mathcal{F}}[l, k] \} - \Im \{ X_{\mathcal{F}}[l, k] \} \end{aligned} \quad (3)$$

It is easy to see that<sup>1</sup>

$$X_S[l, k] = \frac{1}{\sqrt{2}} \left( e^{j \frac{3\pi}{4}} X_{\mathcal{F}}[l, k] + e^{-j \frac{3\pi}{4}} X_{\mathcal{F}}[l, K - k] \right) \quad (4)$$

from which it we deduce

$$X_{\mathcal{F}}[l, k] = \frac{1}{\sqrt{2}} \left( e^{-j \frac{3\pi}{4}} X_S[l, k] + e^{j \frac{3\pi}{4}} X_S[l, K - k] \right) \quad (5)$$

One disadvantage of the SSBT is that the convolution theorem does not hold when employing it (see Appendix A), not even as an approximation. Nonetheless, by converting any result in the SSBT domain to the STFT domain (using Eq. (3)) before utilization, it remains feasible to employ the transform to study of the problem at hand.

### 3. Signal Model and Beamforming

Let there be a device that consists of  $M$  sensors and a loudspeaker (LS) in a reverberant environment. In this setting there also are a desired source, an interfering source, and uncorrelated noise impinging at each sensor, all traveling with a speed  $c$ . For simplicity we assume that all sources are spatially stationary, although condition can be easily removed.

We denote  $y_m[n]$  as the signal at the  $m$ -th sensor, being defined as

$$y_m[n] = h_m[n] * x[n] + g_m[n] * w[n] + e_m[n] * s[n] + r_m[n] \quad (6)$$

in which  $h_m[n]$  is the impulse response between the desired source and the  $m$ -th sensor ( $1 \leq m \leq M$ ), with  $x[n]$  being the desired source's signal; similarly for  $g_m[n]$  and the interfering source  $w[n]$ , and for  $e_m[n]$  and the speaker's signal  $s[n]$ ; and  $r_m[n]$  is the uncorrelated noise.

We let  $m'$  be the reference sensor's index, for simplicity assume  $m' = 1$ , and also  $x_1[n] = h_1[n] * x[n]$  (and similarly for  $v_1[n]$  and  $s_1[n]$ ). We define  $a_m[n]$  as the *relative* impulse response between the desired signal (at the reference sensor) and the  $m$ -th sensor, such that

$$a_m[n] * x_1[n] = h_m[n] * x[n] \quad (7)$$

We similarly define  $b_m[n]$  such that  $b_m[n] * w_1[n] = g_m[n] * w[n]$ , and  $c_m[n]$  from  $e_m[n]$ . Therefore, Eq. (6) becomes

$$y_m[n] = a_m[n] * x_1[n] + b_m[n] * w_1[n] + c_m[n] * s_1[n] + r_m[n] \quad (8)$$

Here, the impulse responses ( $a_m[n]$ ,  $b_m[n]$ ,  $c_m[n]$ ) can be non-causal, depending on the direction of arrival and features of the reverberant environment.

We use a time-frequency transform (such as the STFT or SSBT, as in Section 2) with the convolutive transfer-function (CTF) model [13] to obtain our time-frequency signal model,

$$Y_m[l, k] = A_m[l, k] * X_1[l, k] + B_m[l, k] * W_1[l, k] + C_m[l, k] * S_1[l, k] + R_m[l, k] \quad (9)$$

<sup>1</sup> For the abuse of notation, we let  $X_S[l, K] \equiv X_S[l, 0]$ , and equally for  $X_{\mathcal{F}}[l, K]$ .

where  $Y_m[l, k]$  is the transform of  $y_m[n]$  (resp. all other signals);  $l$  and  $k$  are the window (or decimated-time) and bin indexes, with  $0 \leq k \leq K - 1$ ; and the convolution is in the window-index axis.

Using that  $A_m[l, k]$  is a finite (possibly truncated) response with  $L_A$  windows, then

$$A_m[l, k] * X_1[l, k] = \mathbf{a}_m^T[k] \mathbf{x}_1[l, k] \quad (10)$$

in which

$$\mathbf{a}_m[k] = \left[ A_m[-\Delta, k], \dots, A_m[0, k], \dots, A_m[L_B - \Delta - 1, k] \right]^T \quad (11a)$$

$$\mathbf{x}_1[l, k] = \left[ X_1[l + \Delta, k], \dots, X_1[l, k], \dots, X_1[l - L_B + \Delta + 1, k] \right]^T \quad (11b)$$

and in the same way we define  $\mathbf{b}_m[k]$ ,  $\mathbf{w}_1[l, k]$ ,  $\mathbf{d}_m[k]$  and  $\mathbf{s}_1[l, k]$ . Note that  $\mathbf{a}_m[k]$  and  $\mathbf{b}_m[k]$  don't depend on the index  $l$ , given the spatial stationarity assumption. Also,  $\Delta$  is the number of non-causal windows in the reference sensor necessary to capture the whole signal. With this, Eq. (9) becomes

$$Y_m[l, k] = \mathbf{a}_m^T[k] \mathbf{x}_1[l, k] + \mathbf{b}_m^T[k] \mathbf{w}_1[l, k] + \mathbf{c}_m^T[k] \mathbf{s}_1[l, k] + R_m[l, k] \quad (12)$$

Vectorizing the signals sensor-wise, we finally get

$$\mathbf{y}[l, k] = \mathbf{A}[k] \mathbf{x}_1[l, k] + \mathbf{B}[k] \mathbf{w}_1[l, k] + \mathbf{C}[k] \mathbf{s}_1[l, k] + \mathbf{r}[l, k] \quad (13)$$

where

$$\mathbf{y}[l, k] = \left[ y_1[l, k], \dots, y_M[l, k] \right]^T \quad (14)$$

and similarly for the other variables. In this situation,  $\mathbf{A}[k]$ ,  $\mathbf{B}[k]$  and  $\mathbf{C}[k]$  are  $M \times L_A$ ,  $M \times L_B$  and  $M \times L_C$  matrices respectively;  $\mathbf{x}_1[l, k]$ ,  $\mathbf{w}_1[l, k]$  and  $\mathbf{s}_1[l, k]$  are  $L_A \times 1$ ,  $L_B \times 1$  and  $L_C \times 1$  vectors respectively; and  $\mathbf{y}[l, k]$  and  $\mathbf{r}[l, k]$  are  $M \times 1$  vectors.

### 3.1. Reverb-aware formulation

Let  $l = 0$  be the desired window which we would like to retrieve from the signal  $\mathbf{A}[k] \mathbf{x}_1[l, k]$ . We can write  $\mathbf{A}[k] \mathbf{x}_1[l, k]$  as

$$\mathbf{A}[k] \mathbf{x}_1[l, k] = \mathbf{d}_x[k] X_1[l, k] + \mathbf{q}[l, k] \quad (15)$$

where  $X_1[l, k]$  is the desired speech signal, and  $\mathbf{q}[l, k]$  is an undesired component, uncorrelated to  $X_1[l, k]$ . Through a similar process as exposed in [19] (sec. 7.1.1) (see Appendix B for details), we can deduce that  $\mathbf{d}_x[k]$  can be defined as

$$\mathbf{d}_x[k] = \frac{\sum_i \mathbf{a}_m[k]_i \sum_{n=0}^{K-1-|i\mathcal{O}|} w[n] w[n + |i\mathcal{O}|]}{\sum_{n'=0}^{K-1} w[n']^2} \quad (16)$$

in which  $w[n]$  and  $\mathcal{O}$  are the window-function and the decimation factor used for the time-frequency transform; and  $\mathbf{a}_m[k]_i$  is the  $(\Delta + i)$ -th element of  $\mathbf{a}_m[k]$ .

Note that  $\mathbf{d}_x[k] X_1[l, k]$  also consists of some reverberation, since  $x_1[n] = h_1[n] * x[n]$  is the desired signal at the reference sensor and not at source, therefore being affected by the environment. However, this formulation allows us to better estimate the influence of the neighboring time-frequency windows over the desired signal, given their overlap.

It's important to have in mind the sensor delay and window length. If the time for the signal to travel from the reference to the farthest sensor exceeds the window length (in seconds), multiple windows may represent the desired speech. This isn't a problem if  $\frac{\delta}{c} < \frac{K}{f_s}$ , where  $\delta$  is the biggest reference-to-sensor distance, and  $K$  is the window length.

Using Eq. (15) we define  $\mathbf{v}[l, k]$  as the undesired signal (undesired speech components + speaker signal + interfering source + uncorrelated noise),

$$\mathbf{v}[l, k] = \mathbf{q}[l, k] + \mathbf{B}[k]\mathbf{w}_1[l, k] + \mathbf{C}[k]\mathbf{s}_1[l, k] + \mathbf{r}[l, k] \quad (17)$$

and therefore

$$\mathbf{y}[l, k] = \mathbf{d}_x[k]X_1[l, k] + \mathbf{v}[l, k] \quad (18)$$

We estimate the desired signal at reference  $X_1[l, k]$  as  $Z[l, k]$  through a filter  $\mathbf{f}[l, k]$ , such that

$$\begin{aligned} Z[l, k] &= \mathbf{f}^H[l, k]\mathbf{y}[l, k] \\ &\approx X_1[l, k] \end{aligned} \quad (19)$$

with  $(\cdot)^H$  being the transposed-complex-conjugate operator. Since the source signals' properties can vary over time, so can the filter, adapting to the environment.

In order to get the most minimization on the LS signal, we will use the knowledge of  $\mathbf{C}[k]$  to cancel its windows of most energy. Let  $\mathbf{q}_{l'}[k]$  be the permuted columns of  $\mathbf{C}[k]$  ( $0 \leq l' \leq L_C - 1$ ), such that  $l' < l''$  implies that  $\mathbf{q}_{l'}^H[k]\mathbf{q}_{l''}[k] \geq \mathbf{q}_{l''}^H[k]\mathbf{q}_{l'}[k]$ . We then choose the first  $P < M$  vectors to be nulls of our beamformer,

$$\mathbf{f}^H[l, k]\mathbf{q}_{l'}[k] = 0, \quad 0 \leq l' \leq P - 1 \quad (20)$$

With these constraints, we ensure the erasure of the  $P$  most important windows of  $\mathbf{C}[k]$  from the output signal. These  $P$  constraints, together with the distortionless constraint, give us the  $M \times (P + 1)$  constraint matrix  $\mathbf{P}[k]$ ,

$$\mathbf{f}^H[l, k]\mathbf{P}[k] = \mathbf{i}_{P+1} \quad (21)$$

where  $\mathbf{i}_{P+1} = [1, 0, \dots, 0]$  is a  $(P + 1) \times 1$  vector.

To minimize  $\mathbf{v}[l, k]$  under the constraints from Eq. (21), a Linearly-Constraint Minimum-Variance (LCMV) beamformer will be used, it being defined as

$$\mathbf{f}_{\text{lcmv}}[l, k] = \min_{\mathbf{f}[l, k]} \mathbf{f}[l, k]^H \Phi_{\mathbf{v}}[l, k] \mathbf{f}[l, k] \text{ s.t. } \mathbf{f}^H[l, k]\mathbf{P}[k] = \mathbf{i}_{P+1}^T \quad (22)$$

where  $\Phi_{\mathbf{v}}[l, k]$  is the correlation matrix of the undesired signal  $\mathbf{v}[l, k]$ . The solution to this minimization problem

$$\mathbf{f}_{\text{lcmv}}[l, k] = \Phi_{\mathbf{v}}^{-1}[l, k]\mathbf{P}[k] \left( \mathbf{P}^H[k]\Phi_{\mathbf{v}}^{-1}[l, k]\mathbf{P}[k] \right)^{-1} \mathbf{i}_{P+1} \quad (23)$$

Trivially, the LCMV beamformer requires that  $P + 1 \leq M$ .

### 3.2. Performance metrics

Given the three main goals of the beamformer being the cancellation of the LS signal, the minimization of the overall undesired signal, and the maintenance of the desired signal (through the distortionless constraint), our choice of metrics will reflect these objectives. We define  $S_f[l, k] = \mathbf{f}^H[l, k]\mathbf{s}_1[l, k]$ ,  $X_f[l, k] = \mathbf{f}^H[l, k]\mathbf{x}_1[l, k]$  and  $V_f[l, k] = \mathbf{f}^H[l, k]\mathbf{v}_1[l, k]$  as

the filtered-LS, filtered-desired and filtered-undesired signals, respectively. Unless stated otherwise, the metrics used are broadband. 150  
151

The LS signal's minimization will be measured via the echo-return loss enhancement  $\zeta_s[l]$ , defined as 152  
153

$$\zeta_s[l] = \frac{\sum_k |S_1[l, k]|^2}{\sum_k |S_f[l, k]|^2} \quad (24)$$

The desired signal distortion index  $v[l]$  will be used to assess the distortion on the desired signal, given by 154  
155

$$v[l] = \frac{\sum_k |X_1[l, k] - X_f[l, k]|^2}{\sum_k |X_1[l, k]|^2} \quad (25)$$

Finally, the minimization of the overall undesired signal will be measured using the array gain, such that 156  
157

$$\text{gSNR}[l] = \frac{\sum_k |X_f[l, k]|^2}{\sum_k |V_f[l, k]|^2} \div \frac{\sum_k |X_1[l, k]|^2}{\sum_k |V_1[l, k]|^2} \quad (26)$$

in which  $V_1[l, k]$  is the undesired signal at the reference sensor. 158

#### 4. True-LCMV with the Single-Sideband Transform 159

When carelessly using any of the established methods with the SSBT, the distortionless constraint ensures that the beamformer avoids causing distortion exclusively within the SSBT domain. However, as explained in Section 2 the SSBT beamformer must be carefully constructed to achieve the desired effects, such as the distortionless constraint. 160  
161  
162  
163

We thus propose a framework for the SSBT in which we consider the bins  $k$  and  $K - k$  simultaneously, since from Eq. (5) they both contribute to the  $k$ -th bin in the STFT domain. We define  $\mathbf{w}'[l, k]$  as 164  
165  
166

$$\mathbf{w}'[l, k] = \begin{bmatrix} \mathbf{w}[l, k] \\ \mathbf{w}[l, K - k] \end{bmatrix}_{2M \times 1} \quad (27)$$

from which we define  $\Phi_{\mathbf{w}'}[l, k]$  as its correlation matrix. Under this idea, our filter  $\mathbf{f}'[l, k]$  is a  $2M \times 1$  vector, with the first  $M$  values being for the  $k$ -th bin, and the last  $M$  values for the  $[K - k]$ -th bin. We let the STFT-equivalent filter for the SSBT beamformer  $\mathbf{f}'[l, k]$  be  $\mathbf{f}'_{\mathcal{F}}[l, k]$ , given by 167  
168  
169  
170

$$\mathbf{f}'_{\mathcal{F}}[l, k] = \Lambda \mathbf{f}'[l, k] \quad (28)$$

in which 171

$$\Lambda = \frac{1}{\sqrt{2}} \begin{bmatrix} e^{-j\frac{3\pi}{4}} & 0 & \dots & 0 & e^{j\frac{3\pi}{4}} & 0 & \dots & 0 \\ 0 & e^{-j\frac{3\pi}{4}} & \dots & 0 & 0 & e^{j\frac{3\pi}{4}} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & 0 \\ 0 & 0 & \dots & e^{-j\frac{3\pi}{4}} & 0 & 0 & \dots & e^{j\frac{3\pi}{4}} \end{bmatrix}_{M \times 2M} \quad (29)$$

From Eqs. (21) and (28) the constraint matrix within the SSBT domain becomes 172

$$\mathbf{f}'^H[l, k] \hat{\mathbf{P}}[k] = \mathbf{i}_{P+1}^T \quad (30a)$$

$$\hat{\mathbf{P}}[k] = \mathbf{\Lambda}^H \mathbf{P}_{\mathcal{F}}[k] \quad (30b)$$

where  $\mathbf{P}_{\mathcal{F}}[l, k]$  is the constraint matrix within the STFT domain.

In this scheme, our minimization problem becomes

$$\mathbf{f}'_{\text{lcmv}}[l, k] = \min_{\mathbf{f}'[l, k]} \mathbf{f}'^H[l, k] \mathbf{\Phi}_{\mathbf{w}'}[l, k] \mathbf{f}'[l, k] \text{ s.t. } \mathbf{f}'^H[l, k] \hat{\mathbf{P}}[k] = \mathbf{i}_{P+1} \quad (31)$$

Although  $\mathbf{\Phi}_{\mathbf{w}'}[l, k]$  is a matrix with real entries,  $\hat{\mathbf{P}}[k]$  is complex-valued, and thus is the solution to Eq. (31), contradicting the purpose of utilizing the SSBT.

#### 4.1. Real-valued true-LCMV beamformer with SSBT

To ensure the desired behavior of  $\mathbf{f}'[l, k]$  being real-valued, an additional constraint is necessary. By forcing  $\mathbf{f}'[l, k]$  to have real entries, from Eq. (30a) we trivially have that

$$\mathbf{f}'^T[l, k] \Re\{\hat{\mathbf{P}}[k]\} = \mathbf{i}_{P+1}^T \quad (32a)$$

$$\mathbf{f}'^T[l, k] \Im\{\hat{\mathbf{P}}[k]\} = \mathbf{0}_{P+1}^T \quad (32b)$$

which can be put in matricial form as  $\mathbf{f}'^T[l, k] \hat{\mathbf{P}}'[k] = \mathbf{i}_{2(P+1)}^T$ , with

$$\hat{\mathbf{P}}'[k] = \begin{bmatrix} \Re\{\hat{\mathbf{P}}[k]\} & \Im\{\hat{\mathbf{P}}[k]\} \end{bmatrix}_{2M \times 2(P+1)} \quad (33a)$$

$$\mathbf{i}_{2(P+1)} = [1, 0, \dots, 0]^T \quad (33b)$$

Therefore, the minimization problem from Eq. (31) becomes

$$\mathbf{f}'_{\text{lcmv}}[l, k] = \min_{\mathbf{f}'[l, k]} \mathbf{f}'^T[l, k] \mathbf{\Phi}_{\mathbf{w}'}[l, k] \mathbf{f}'[l, k] \text{ s.t. } \mathbf{f}'^T[l, k] \hat{\mathbf{P}}'[k] = \mathbf{i}_{2(P+1)}^T \quad (34)$$

whose solution is

$$\mathbf{f}'_{\text{lcmv}}[l, k] = \mathbf{\Phi}_{\mathbf{w}'}^{-1}[l, k] \hat{\mathbf{P}}' x[k] \left( \hat{\mathbf{P}}' x^T[k] \mathbf{\Phi}_{\mathbf{w}'}^{-1}[k] \hat{\mathbf{P}}' x[k] \right)^{-1} \mathbf{i}_{2(P+1)} \quad (35)$$

Using Eq. (28), we can obtain the desired beamformer  $\mathbf{f}'_{\mathcal{F}, \text{lcmv}}[l, k]$ , transformed to the STFT domain.

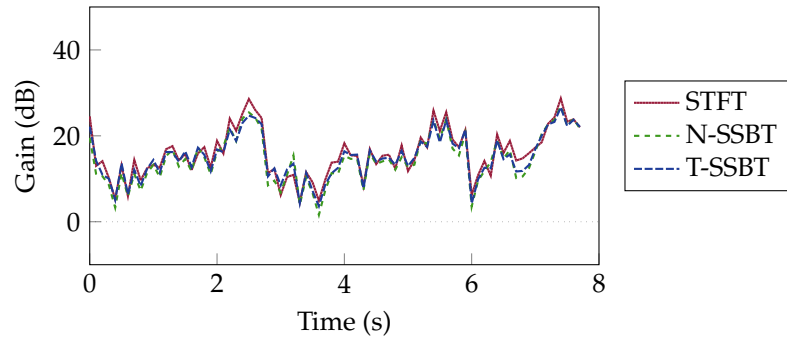
## 5. Simulations

In the simulations<sup>2</sup>, we employ a sampling frequency of 16kHz. The sensor array consists of a uniform linear array with 10 sensors spaced at 2cm. Room impulse responses were generated using Habets' RIR generator [20], and signals were selected from the SMARD [21] and LINSE [22] databases.

The room's dimensions are 4m × 6m × 3m (width × length × height), with a reverberation time of 0.3s. The desired source is located at (2m, 1m, 1m), it being a male voice (SMARD, 50\_male\_speech\_english\_ch8\_OmniPower4296.flac). The device composed of the loudspeaker + sensors is centered at (2m, 2m, 1m). It is comprised of 10 radially-symmetric sensors located 8cm away from the center, all omnidirectional and of flat frequency response. In the center is the loudspeaker source, whose signal is a music noise (SMARD database, 69\_abba\_ch8\_OmniPower4296.flac). The interfering source, simulating an open door, is located simultaneously at (0.5m, 5m, [0.3 : 0.3 : 2.7]m), with a babble sound signal (LINSE database, babble.mat). The noise signal is white Gaussian noise (SMARD database, wgn\_48kHz\_ch8\_OmniPower4296.flac). All signals were resampled to the desired frequency.

At the reference sensor, the SNR for the loudspeaker's signal, interfering signal and noise are, respectively, of −15dB, 10dB and 30dB. The beamformers were calculated every 200 windows, and used up to the previous 1000 samples to estimate the correlation matrices.

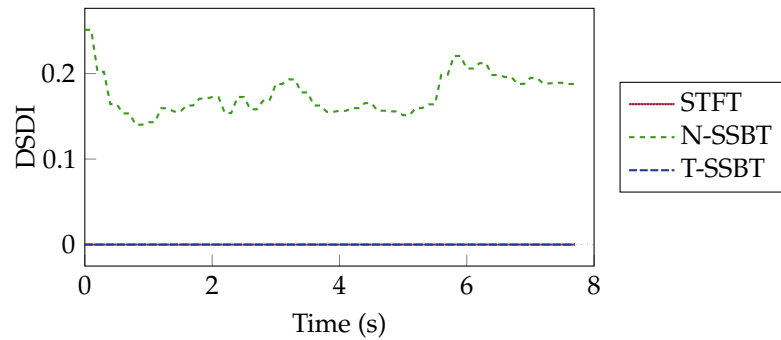
We compare filters obtained through the STFT and SSBT transforms. N-SSBT uses Eq. (23) to (naively) calculate the SSBT beamformer, and T-SSBT will denote the beamformer obtained via the true-distortionless MVDR from Section 4. Performance analysis is conducted via the STFT domain, with the SSBT beamformers being converted into it. In line plots, STFT is presented in red, N-SSBT in green, and T-SSBT in blue. The output metrics were averaged over 20 windows, to facilitate visualization.



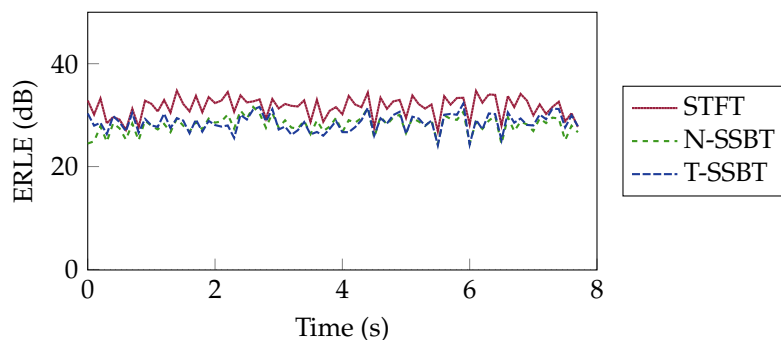
**Figure 1.** Per-window broadband gain for  $K = 32$ .

<sup>2</sup> Code is available at <https://github.com/VCurtarelli/py-ssb-ctf-bf>.





**Figure 2.** Per-window broadband DSDI gain for  $K = 32$ .



**Figure 3.** Per-window broadband ERLE gain for  $K = 32$ .

From Fig. 1, all beamformers achieved a similar gain in SNR over time, with some fluctuations but no distinguishable advantage. From Fig. 2, both the STFT and T-SSBT beamformers achieved the desired null distortion on the desired signal, while the N-SSBT one caused some minimal distortion, although present.

In Fig. 3, we see that the STFT beamformer clearly outperformed both other beamformers in terms of ERLE, which would mean a better reduction of the loudspeaker's signal, the main objective.

**Author Contributions:** Conceptualization, I. Cohen and V. Curtarelli; Methodology, V. Curtarelli; Software, V. Curtarelli; Writing—original draft: V. Curtarelli; Writing—review and editing, I. Cohen and V. Curtarelli; Supervision, V. Curtarelli. All authors have read and agreed to the published version of the manuscript.

**Funding:** This research was supported by the Pazy Research Foundation, and the Israel Science Foundation (grant no. 1449/23).

**Data Availability Statement:** The source-code for the simulations developed here is available at <https://github.com/VCurtarelli/py-ssb-ctf-bf>.

**Conflicts of Interest:** The authors declare no conflict of interest.

## Abbreviations

The following abbreviations are used in this manuscript:

CTF	Convolutional Transfer Function
DSRF	Desired Signal Reduction Factor
LCMV	Linearly-Constrained Minimum-Variance
MVDR	Minimum-Variance Distortionless-Response
MTF	Multiplicative Transfer Function
SNR	Signal-to-Noise Ratio
SSBT	Single-Sideband Transform
STFT	Short-Time Fourier Transform

## Appendix A. SSBT Convolution

Let  $x[n]$  be a time domain signal, with  $X_{\mathcal{F}}[l, k]$  being its STFT equivalent, and  $X_{\mathcal{S}}[l, k]$  its SSBT equivalent. We here assume that both the STFT and the SSBT have  $K$  frequency bins.  $X_{\mathcal{S}}[l, k]$  can be obtained using  $X_{\mathcal{F}}[l, k]$ , through

$$X_{\mathcal{S}}[l, k] = -X_{\mathcal{F}}^{\Re}[l, k] - X_{\mathcal{F}}^{\Im}[l, k] \quad (\text{A1})$$

in which  $(\cdot)^{\Re}$  and  $(\cdot)^{\Im}$  represent the real and imaginary components of their argument, respectively.

It is easy to see that

$$X_{\mathcal{F}}[l, k] = \frac{1}{\sqrt{2}} \left( e^{-j\frac{3\pi}{4}} X_{\mathcal{S}}[l, k] + e^{j\frac{3\pi}{4}} X_{\mathcal{S}}[l, K - k] \right) \quad (\text{A2})$$

As stated before, in this formulation we abuse the notation by letting  $X_{\mathcal{S}}[l, K] = X_{\mathcal{S}}[l, 0]$  to simplify the mathematical operations.

Now, let there also be  $h[n]$ ,  $H_{\mathcal{F}}[k]$  and  $H_{\mathcal{S}}[k]$ , with the same assumptions as before. We define  $Y_{\mathcal{F}}[l, k]$  and  $Y_{\mathcal{S}}[l, k]$  as the output of an LTI system with impulse response  $h[n]$ , such that

$$Y_{\mathcal{F}}[l, k] = H_{\mathcal{F}}[k] X_{\mathcal{F}}[l, k] \quad (\text{A3a})$$

$$Y_{\mathcal{S}}[l, k] = H_{\mathcal{S}}[k] X_{\mathcal{S}}[l, k] \quad (\text{A3b})$$

We will assume that the MTF model [13] correctly models the convolution here. This was used instead of the CTF for simplicity, as these derivations would work exactly the same for the CTF, but with window-wise summations as well, which would pollute the notation.

Applying Eq. (A1) in Eq. (A3b), and knowing that  $X_{\mathcal{F}}[l, k] = X_{\mathcal{F}}^*[l, K - k]$  (same for  $H_{\mathcal{F}}[k]$ ), with  $(\cdot)^*$  representing the complex-conjugate; we get that

$$\begin{aligned} Y_{\mathcal{S}}[l, k] &= X_{\mathcal{F}}^{\Re}[l, k] H_{\mathcal{F}}^{\Re}[l, k] + X_{\mathcal{F}}^{\Re}[l, k] H_{\mathcal{F}}^{\Im}[l, k] \\ &\quad + X_{\mathcal{F}}^{\Im}[l, k] H_{\mathcal{F}}^{\Re}[l, k] + X_{\mathcal{F}}^{\Im}[l, k] H_{\mathcal{F}}^{\Im}[l, k] \\ Y_{\mathcal{S}}[l, K - k] &= X_{\mathcal{F}}^{\Re}[l, k] H_{\mathcal{F}}^{\Re}[l, k] - X_{\mathcal{F}}^{\Re}[l, k] H_{\mathcal{F}}^{\Im}[l, k] \\ &\quad - X_{\mathcal{F}}^{\Im}[l, k] H_{\mathcal{F}}^{\Re}[l, k] + X_{\mathcal{F}}^{\Im}[l, k] H_{\mathcal{F}}^{\Im}[l, k] \end{aligned} \quad (\text{A4})$$

Passing this through Eq. (A2),

$$\begin{aligned} Y'_{\mathcal{F}}[l, k] &= -X_{\mathcal{F}}^{\Re}[l, k] H_{\mathcal{F}}^{\Re}[l, k] + j X_{\mathcal{F}}^{\Re}[l, k] H_{\mathcal{F}}^{\Re}[l, k] \\ &\quad + j X_{\mathcal{F}}^{\Im}[l, k] H_{\mathcal{F}}^{\Re}[l, k] - X_{\mathcal{F}}^{\Im}[l, k] H_{\mathcal{F}}^{\Im}[l, k] \end{aligned} \quad (\text{A5})$$

where  $Y'_{\mathcal{F}}[l, k]$  is the STFT-equivalent of  $Y_{\mathcal{S}}[l, k]$ .

Expanding Eq. (A3a) in terms of real and imaginary components,

$$\begin{aligned} Y_{\mathcal{F}}[l, k] &= X_{\mathcal{F}}^{\Re}[l, k]H_{\mathcal{F}}^{\Re}[l, k] + jX_{\mathcal{F}}^{\Re}[l, k]H_{\mathcal{F}}^{\Im}[l, k] \\ &\quad + jX_{\mathcal{F}}^{\Im}[l, k]H_{\mathcal{F}}^{\Re}[l, k] - X_{\mathcal{F}}^{\Im}[l, k]H_{\mathcal{F}}^{\Im}[l, k] \end{aligned} \quad (\text{A6})$$

Comparing Eq. (A5) and Eq. (A6), trivially  $Y'_{\mathcal{F}}[l, k] \neq Y_{\mathcal{F}}[l, k]$ . This proves that the SSBT doesn't appropriately models the convolution, and therefore the convolution theorem doesn't hold when applying this transform.

## Appendix B. Correct separation of desired signal

Let  $X_m[l, k]$  be such that

$$\begin{aligned} X_m[l, k] &= A_m[l, k] * X_1[l, k] \\ &= \mathbf{a}_m^{\top}[k] \mathbf{x}_1[l, k] \end{aligned} \quad (\text{A7})$$

as in Eqs. (9) and (10). We can separate  $X_m[l, k]$  as

$$X_m[l, k] = d_m[l, k]X_1[l, k] + X'[l, k] \quad (\text{A8})$$

where  $d_m[l, k]$  is the steering vector for the desired speech portion  $X_1[l, k]$ , and  $X'[l, k]$  is the undesired speech component, in such a way that  $X_1[l, k]$  and  $X'[l, k]$  are uncorrelated.

This seems trivial in a first glance, by adopting  $d_m[k]$  as the 0-th element of  $\mathbf{a}_m[k]$ , and  $X'[l, k]$  as the rest of the summation. However, this doesn't take into account that  $X[l, k]$  and  $X[l', k]$  may be correlated if  $|l - l'| \leq S$ , where  $S = \lfloor (K-1)/\mathcal{O} \rfloor$ . That is, if  $l$ -th and  $l'$ -th transform window share samples, then there is some correlation between them.

We define

$$\begin{aligned} d_m[l, k] &= \frac{\mathbb{E}\{X_1^*[l, k]X_m[l, k]\}}{\mathbb{E}\{|X_1[l, k]|^2\}} \\ &= \frac{\sum_i A_m[i, k]\mathbb{E}\{X_1^*[l, k]X_1[l - i, k]\}}{\mathbb{E}\{|X_1[l, k]|^2\}} \end{aligned} \quad (\text{A9})$$

Focusing on each expectation in the numerator,

$$E_i = \mathbb{E}\{X_1^*[l, k]X_1[l - i, k]\} \quad (\text{A10})$$

Through the definition of the STFT,

$$\begin{aligned} E_i &= \mathbb{E}\left\{\sum_{n=0}^{K-1} w(n)x_1(n + l\mathcal{O})e^{j2\pi\frac{k}{K}(n+l\mathcal{O})} \sum_{v=0}^{K-1} w(v)x_1(v + (l - i)\mathcal{O})e^{-j2\pi\frac{k}{K}(v+(l-i)\mathcal{O})}\right\} \\ &= \sum_{n=0}^{K-1} \sum_{v=0}^{K-1} w(n)w(v)e^{-j2\pi\frac{k}{K}(v-n+i\mathcal{O})} \mathbb{E}\{x_1(n + l\mathcal{O})x_1(v + (l - i)\mathcal{O})\} \end{aligned} \quad (\text{A11})$$

Using the substitutions  $\tilde{n} = n + l\mathcal{O}$  and  $\tilde{v} = v + (l - i)\mathcal{O}$ ,

$$E_i = \sum_{\tilde{v}=l\mathcal{O}}^{K-1+l\mathcal{O}} \sum_{\tilde{n}=(l-i)\mathcal{O}}^{K-1+(l-i)\mathcal{O}} w(\tilde{n} - l\mathcal{O})w(\tilde{v} - (l - i)\mathcal{O})\mathbb{E}\{x_1(\tilde{n})x_1(\tilde{v})\}e^{-j2\pi\frac{k}{K}(\tilde{v}-\tilde{n})} \quad (\text{A12})$$

We now assume  $x_1(n)$  is the result of a zero-mean white process, and therefore  $x_1(\tilde{n})$  and  $x_1(\tilde{v})$  are independent if  $\tilde{n} \neq \tilde{v}$  (which isn't true, but will allow us to continue the derivations). Then

$$E_i = \sum_{\tilde{n}} w(\tilde{n} - l\mathcal{O})w(\tilde{n} - (l - i)\mathcal{O})E\{x_1(\tilde{n})^2\} \quad (\text{A13})$$

Rolling back the substitution  $n = \tilde{n} - l\mathcal{O}$ ,

$$E_i = \sum_{n=0}^{K-1-|i\mathcal{O}|} w(n)w(n + i\mathcal{O})E\{x_1(n + l\mathcal{O})^2\} \quad (\text{A14})$$

where the summation takes into account that both windows are finite. Now we at last assume that  $E\{x_1(n + l\mathcal{O})^2\} \approx E\{x_1(l\mathcal{O})^2\}$  for small values of  $n$  (such as within at most one window), such that

$$E_i = \phi_{x_1}(l\mathcal{O}) \sum_{n=0}^{K-1-|i\mathcal{O}|} w(n)w(n + i\mathcal{O}) \quad (\text{A15})$$

Going back to Eq. (A9), and applying the derivations also to the numerator (with  $i = 0$ ),

$$\begin{aligned} d_m[l, k] &= \frac{\sum_i A_m[i, k] \phi_{x_1}(l\mathcal{O}) \sum_{n=0}^{K-1-|i\mathcal{O}|} w(n)w(n + |i\mathcal{O}|)}{\phi_{x_1}(l\mathcal{O}) \sum_{n'=0}^{K-1} w(n')^2} \\ &= \frac{\sum_i \sum_{n=0}^{K-1-|i\mathcal{O}|} A_m[i, k] w(n)w(n + |i\mathcal{O}|)}{\sum_{n'=0}^{K-1} w(n')^2} \end{aligned} \quad (\text{A16})$$

This way, we see that  $d_m[l, k] \equiv d_m[k]$  since it doesn't depend on  $l$ . However, since for the summation over  $n$  we need that  $K - 1 - |i\mathcal{O}| \geq 0$ , this means that  $-S \leq i \leq S$ , and thus our desired speech signal  $X_1[l, k]$  depends on up to the previous (or next)  $S = \lfloor (K-1)/\mathcal{O} \rfloor$  windows.

## References

1. Lobato, W.; Costa, M.H. Worst-Case-Optimization Robust-MVDR Beamformer for Stereo Noise Reduction in Hearing Aids. *IEEE/ACM Transactions on Audio, Speech, and Language Processing* **2020**, *28*, 2224–2237. <https://doi.org/10.1109/TASLP.2020.3009831>.
2. Chen, J.; Kung Yao, R.; Hudson, R. Source localization and beamforming. *IEEE Signal Processing Magazine* **2002**, *19*, 30–39. <https://doi.org/10.1109/79.985676>.
3. Lu, J.Y.; Zou, H.; Greenleaf, J.F. Biomedical ultrasound beam forming. *Ultrasound in Medicine & Biology* **1994**, *20*, 403–428. [https://doi.org/10.1016/0301-5629\(94\)90097-3](https://doi.org/10.1016/0301-5629(94)90097-3).
4. Nguyen, N.Q.; Prager, R.W. Minimum Variance Approaches to Ultrasound Pixel-Based Beamforming. *IEEE Transactions on Medical Imaging* **2017**, *36*, 374–384. <https://doi.org/10.1109/TMI.2016.2609889>.
5. Benesty, J.; Cohen, I.; Chen, J. *Fundamentals of signal enhancement and array signal processing*; John Wiley & Sons: Hoboken, NJ, 2017.
6. Kıymık, M.; Güler, İ.; Dizibüyük, A.; Akın, M. Comparison of STFT and wavelet transform methods in determining epileptic seizure activity in EEG signals for real-time application. *Computers in Biology and Medicine* **2005**, *35*, 603–616. <https://doi.org/10.1016/j.compbiomed.2004.05.001>.
7. Pan, C.; Chen, J.; Shi, G.; Benesty, J. On microphone array beamforming and insights into the underlying signal models in the short-time-Fourier-transform domain. *The Journal of the Acoustical Society of America* **2021**, *149*, 660–672. <https://doi.org/10.1121/10.0003335>.
8. Chen, W.; Huang, X. Wavelet-Based Beamforming for High-Speed Rotating Acoustic Source. *IEEE Access* **2018**, *6*, 10231–10239. <https://doi.org/10.1109/ACCESS.2018.2795538>.

9. Yang, Y.; Peng, Z.K.; Dong, X.J.; Zhang, W.M.; Meng, G. General Parameterized Time-Frequency Transform. *IEEE Transactions on Signal Processing* **2014**, *62*, 2751–2764. <https://doi.org/10.1109/TSP.2014.2314061>.
10. Almeida, L. The fractional Fourier transform and time-frequency representations. *IEEE Transactions on Signal Processing* **1994**, *42*, 3084–3091. <https://doi.org/10.1109/78.330368>.
11. Crochiere, R.E.; Rabiner, L.R. *Multirate digital signal processing*; Prentice-Hall signal processing series, Prentice-Hall: Englewood Cliffs, N.J, 1983.
12. Oyerman, A. Speech Dereverberation in the Time-Frequency Domain. Master's thesis, Technion - Israel Institute of Technology, Haifa, Israel, 2012.
13. Talmon, R.; Cohen, I.; Gannot, S. Relative Transfer Function Identification Using Convolutional Transfer Function Approximation. *IEEE Transactions on Audio, Speech, and Language Processing* **2009**, *17*, 546–555. <https://doi.org/10.1109/TASL.2008.2009576>.
14. Kumatani, K.; McDonough, J.; Schacht, S.; Klakow, D.; Garner, P.N.; Li, W. Filter bank design based on minimization of individual aliasing terms for minimum mutual information subband adaptive beamforming. In Proceedings of the 2008 IEEE International Conference on Acoustics, Speech and Signal Processing, Las Vegas, NV, March 2008; pp. 1609–1612. <https://doi.org/10.1109/ICASSP.2008.4517933>.
15. Gopinath, R.; Burrus, C. A tutorial overview of filter banks, wavelets and interrelations. In Proceedings of the 1993 IEEE International Symposium on Circuits and Systems, Chicago, IL, USA, 1993; pp. 104–107. <https://doi.org/10.1109/ISCAS.1993.393668>.
16. Capon, J. High-resolution frequency-wavenumber spectrum analysis. *Proceedings of the IEEE* **1969**, *57*, 1408–1418. <https://doi.org/10.1109/PROC.1969.7278>.
17. Erdogan, H.; Hershey, J.R.; Watanabe, S.; Mandel, M.I.; Roux, J.L. Improved MVDR Beamforming Using Single-Channel Mask Prediction Networks. In Proceedings of the Interspeech 2016. ISCA, September 2016, pp. 1981–1985. <https://doi.org/10.21437/Interspeech.2016-552>.
18. DeMuth, G. Frequency domain beamforming techniques. In Proceedings of the ICASSP '77. IEEE International Conference on Acoustics, Speech, and Signal Processing, Hartford, CT, USA, 1977; Vol. 2, pp. 713–715. <https://doi.org/10.1109/ICASSP.1977.1170316>.
19. Bai, M.R.; Ih, J.G.; Benesty, J. *Acoustic Array Systems: Theory, Implementation, and Application*, 1 ed.; Wiley, 2013. <https://doi.org/10.1002/9780470827253>.
20. Habets, E. RIR Generator, 2020.
21. Nielsen, J.K.; Jensen, J.R.; Jensen, S.H.; Christensen, M.G. The Single- and Multichannel Audio Recordings Database (SMARD). In Proceedings of the Int. Workshop Acoustic Signal Enhancement, Sep. 2014.
22. Johnson, D.H. Signal Processing Information Database, 2013.

**Disclaimer/Publisher's Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.