# Cyclic Multichannel Wiener Filter for Acoustic Beamforming

*Giovanni Bologni[♯], Richard Heusdens[♯ ♭] and Richard C. Hendriks[♯]*

[♯] Delft University of Technology, Delft, the Netherlands
[♭] Netherlands Defence Academy, Den Helder, the Netherlands

*Abstract*—**Acoustic beamforming models typically assume wide-sense stationarity of speech signals within short time frames. However, voiced speech is better modeled as a cyclostationary (CS) process, a random process whose mean and autocorrelation are $T_1$-periodic, where $\alpha_1 = 1/T_1$ corresponds to the fundamental frequency of vowels. Higher harmonic frequencies are found at integer multiples of the fundamental. This work introduces a cyclic multichannel Wiener filter (cMWF) for speech enhancement derived from a cyclostationary model. This beamformer exploits spectral correlation across the harmonic frequencies of the signal to further reduce the mean-squared error (MSE) between the target and the processed input. The proposed cMWF is optimal in the MSE sense and reduces to the MWF when the target is wide-sense stationary. Experiments on simulated data demonstrate considerable improvements in scale-invariant signal-to-distortion ratio (SI-SDR) on synthetic data but also indicate high sensitivity to the accuracy of the estimated fundamental frequency $\alpha_1$, which limits effectiveness on real data.**

## 1. INTRODUCTION

A noticeable trait of speech is non-stationarity. To address non-stationarity, recordings are often divided into short segments, which are then modeled as realizations of wide-sense stationary (WSS) processes in applications such as dereverberation and beamforming [1]–[4]. However, because of the nearly periodic pressure waves generated by the movement of the vocal folds, voiced speech segments do not behave like WSS processes. Recently, it has been shown that voiced speech can better be modeled as a (wide-sense) *cyclostationary* (CS) process [5]. We will refer to wide-sense cyclostationary processes simply as cyclostationary (CS) in the following. CS processes describe random signals with first- and second-order moments that vary with frequency $\alpha_1$ [6]–[9]. A defining characteristic of CS processes is to exhibit statistical correlation across frequencies, in direct contrast to the WSS assumption, which assumes signals to be uncorrelated over frequency [10], [11]. This distinction is particularly relevant for signals such as voiced speech, where harmonic components at integer multiples of the fundamental frequency $\alpha_1$ occur simultaneously [12, Ch. 8]. In this context, $\alpha_1$ corresponds directly to the fundamental frequency of vowels.

Adaptive filters typically take advantage of temporal correlations in the signals, whereas spatial filters, known as beamformers, exploit spatial correlations. This work is grounded in the theory of FREquency-SHifted (FRESH) filtering, which reconstructs CS signals corrupted by noise exploiting the statistical correlation between cyclic frequencies [10]. In other words, FRESH filtering leverages spectral correlations at harmonic frequencies to improve reconstruction accuracy. FRESH filtering can be extended to multichannel acoustic recordings to exploit spatial and spectral correlations jointly [13].

In this paper, we propose a novel beamformer, the cyclic multichannel Wiener filter (cMWF). Similar to the MWF, our beamformer minimizes the mean-squared error (MSE) between the target and the filter output in the short-time Fourier transform (STFT) domain. Unlike the MWF, both spatial and spectral correlation of the desired signal are considered when deriving the optimal weights. The unknown target spectral-spatial covariance matrix is estimated using a generalized

eigenvalue decomposition followed by a low-rank approximation. Our experiments show that the cMWF is particularly advantageous in low-SNR contexts, and it enjoys significant scale-invariant signal-to-distortion ratio (SI-SDR) gains if the signal to reconstruct is indeed cyclostationary at the frequencies of interest. Moreover, the cMWF reduces to the MWF if the signals under analysis are WSS. However, the proposed approach is highly sensitive to errors in the estimation of the fundamental frequency of the target signal, which limits the applicability of the cMWF in speech processing. A Python implementation of all algorithms is available [14].

## 2. BACKGROUND

Let us begin by introducing some theory of CS processes. We will denote random variables by capitals and the corresponding realizations by small letters. Matrices are denoted by bold capitals and vectors by bold small letters. The tilde denotes frequency domain variables. A real-valued discrete-time random process $\{X(n), n \in \mathbb{Z}\}$ is *cyclostationary* (CS) in the wide sense if both its mean and covariance function are periodic with some period $P$:

$$\mu_X(n) = \mu_X(n+P), \quad r_X(n,\tau) = r_X(n+P,\tau), \ \forall n,\tau \in \mathbb{Z}. \quad (1)$$

As the mean and the covariance of a CS process are periodic in $n$ with period $P$, they accept a Fourier series expansion over a set of *cyclic frequencies* $\mathcal{A} = \{\alpha_p : 2\pi p/P, \ p = 0, \dots, P-1\}$. The covariance can thus be expressed as $r_X(n,\tau) = \sum_{\alpha_p \in \mathcal{A}} c_X(\alpha_p, \tau) \exp(j\alpha_p n)$, where the Fourier coefficients, called *cyclic correlations*, are given by $c_X(\alpha_p, \tau) = P^{-1} \sum_{n=0}^{P-1} r_X(n,\tau) \exp(-j\alpha_p n)$. Now, suppose $c_X(\alpha_p, \tau)$ is absolutely summable w. r. t. $\tau$ for all $n$ in $\mathbb{Z}$. By applying a discrete-time Fourier transform to $c_X(\alpha_p, \tau)$, we get a function $S_X(\alpha_p, \omega)$ of two frequency variables, a *cyclic* frequency $\alpha_p$ and a *spectral* frequency $\omega$: $S_X(\alpha_p, \omega) = \sum_{\tau=-\infty}^{\infty} c_X(\alpha_p, \tau) \exp(-j\omega\tau)$ [7]. The quantity $S_X(\alpha_p, \omega)$ is known as *spectral correlation density*, or cyclic spectrum, as for finite-length processes it is also given by:

$$S_X(\alpha_p, \omega) = \mathbb{E}[\tilde{X}_N(\omega)\tilde{X}_N^*(\omega - \alpha_p)], \quad (2)$$

where $\tilde{X}_N(\omega) = \sum_{n=0}^{N-1} X(n) \exp(-j\omega n)$ is the $N$-point Fourier transform of $\{X(n)\}$. The spectral correlation density (SCD) boils down to the conventional power spectral density (PSD) for $p = 0$. A key property of CS processes is to exhibit inter-frequency correlations. In fact, $\tilde{X}_N(\omega_1)$ is correlated with $\tilde{X}_N(\omega_2)$ for $|\omega_1 - \omega_2| = \alpha_p, \ \forall \alpha_p \in \mathcal{A} \setminus \{0\}$. Intuitively, if we measure $\tilde{x}_N(\omega_1)$, we know something about $\tilde{x}_N(\omega_2)$. In contrast, the spectral components of WSS processes are asymptotically uncorrelated. For example, for white Gaussian noise, we have that $S_X(\alpha_p, \omega) = 0$ for all $\alpha_p \neq 0$. Notice that all quantities in this section are defined for a single random process, but generalizing the notions to the cross-statistics between multiple processes is straightforward.

### 2.1. Estimation of the cyclic spectrum

The proposed beamformer, which will be introduced in Section 3, requires knowledge of the cyclic spectrum $S_X(\alpha_p, \omega)$. However, the definition of the cyclic spectrum in (2) involves an ensemble

arXiv:2507.10159v1 [eess.AS] 14 Jul 2025

expectation. To estimate the cyclic spectrum of a CS process, one can use the *time-averaged cyclic periodogram* (ACP) algorithm [15]. Essentially, the ACP replaces the expectation with a time average and coincides with Welch's PSD estimator for $p = 0$ [16]. The ACP estimator has the desirable property to produce consistent estimates of the SCD even from a single record or realization of the signal. Other methods for SCD estimation may offer faster computations if knowledge of the SCD at all spectral and cyclic frequencies is required [17]–[19].

Let $\{X(n), n \in \mathbb{Z}\}$ and $\{Y(n), n \in \mathbb{Z}\}$ be random processes sampled with sampling frequency $f_s$. The processes $\{X_N(n)\}$ and $\{Y_N(n)\}$ equal $\{X(n)\}$ and $\{Y(n)\}$ for $n \in \{0, \ldots, N-1\}$ and are zero otherwise. Processing these signals in the STFT domain, where the window length $K$ equals the DFT points and the block shift is $R$, yields a total of $L = \lceil 1 + (N-K)/R \rceil$ frames. Notice that the spectral resolution is determined by the length $K$ of the DFT analysis window, with $\Delta \omega \approx f_s/K$ [Hz]. In contrast, the cyclic frequencies $\alpha_p$ are sampled on a finer grid. Their resolution depends on the total length of the signal, giving $\Delta \alpha \approx f_s/(LR)$ [Hz] [15]. The frequency shifted components $\tilde{X}(\omega_k - \alpha_p)$ are not $1/K$ separated. Instead, the frequency translation at the right-hand side of (2) is achieved by first modulating in the time domain with cyclic frequency $\alpha_p$ and then transforming to the frequency domain, which takes advantage of the modulation property of the DFT:

$$\tilde{X}(\omega_k - \alpha_p) \xleftrightarrow{\mathcal{F}} X(n)e^{j\alpha_p n}. \tag{3}$$

The modulated signal in the time domain and its STFT counterpart are given by:

$$X_N^{(\alpha_c)}(n) = X_N(n)e^{jn\alpha_p}, \tag{4a}$$

$$\tilde{X}(\omega_k - \alpha_p, \ell) = \sum_{n=0}^{N-1} X_N^{(\alpha_c)}(n + \ell R)w(n)e^{-jn\omega_k}, \tag{4b}$$

where $\ell$ is the time-frame index and $w(n)$ represents a window function of length $N$. The ACP estimate at cyclic frequency $\alpha_p$ and spectral frequency $\omega_k$ is then given by:

$$\hat{S}_{YX}(\alpha_p, \omega_k) = \frac{1}{L} \sum_{\ell=0}^{L-1} \tilde{Y}(\omega_k, \ell)\tilde{X}^*(\omega_k - \alpha_p, \ell). \tag{5}$$

Additional details on implementing the ACP estimator on natural data are discussed in Section 3.2.

### 2.2. Narrowband beamforming

Let us introduce the signal model and briefly review the beamforming theory. The general goal of beamforming is to estimate the target signal as a linear combination of the noisy inputs. Let $\boldsymbol{x}(\omega_k) = [\tilde{X}_0(\omega_k) \ \ldots \ \tilde{X}_{M-1}(\omega_k)]^T \in \mathbb{C}^M$ denote noisy and reverberant measurements from a microphone array with $M$ elements in the STFT domain:

$$\boldsymbol{x}(\omega_k) = \tilde{S}(\omega_k)\,\boldsymbol{a}(\omega_k) + \boldsymbol{v}(\omega_k) = \boldsymbol{d}(\omega_k) + \boldsymbol{v}(\omega_k), \tag{6}$$

where $\boldsymbol{a}(\omega_k) = [1 \ a_1(\omega_k) \ \ldots \ a_{M-1}(\omega_k)]^T$ is the relative transfer function (RTF) between a reference sensor, the first one in this case, and the remaining sensors, $\tilde{S}(\omega_k)$ is the target signal at the reference microphone, and $\boldsymbol{v}(\omega_k)$ is a noise term. Following the multiplicative transfer function approximation, the late reverberation component is neglected [20]. The MWF is a well-known beamformer that minimizes the MSE between the (unknown) target and the input signals [21]. To avoid that the norm of the weights becomes too large in presence

of estimation errors, an L2 regularization term can be added to the reconstruction loss, yielding the so-called *robust MWF* [22]:

$$\min_{\boldsymbol{w}(\omega_k)} \quad \mathbb{E}[|\tilde{S}(\omega_k) - \boldsymbol{w}^H(\omega_k)\boldsymbol{x}(\omega_k)|^2] + \lambda\|\boldsymbol{w}\|_2^2, \tag{7}$$

where $\lambda$ is a hyperparameter that balances the contribution of the two loss terms. The solution to (7) is given by

$$\boldsymbol{w}_{\mathrm{MWF}}(\omega_k) = (\boldsymbol{R}_x(\omega_k) + \lambda\boldsymbol{I})^{-1}\sigma_s^2(\omega_k)\boldsymbol{a}(\omega_k), \tag{8}$$

where $\boldsymbol{R}_x(\omega_k) = \mathbb{E}[\boldsymbol{x}(\omega_k)\boldsymbol{x}^H(\omega_k)]$ is the noisy covariance matrix and $\sigma_s^2(\omega_k) = \mathbb{E}[|\tilde{S}(\omega_k)|^2]$ is the variance of the target signal. A possible choice of $\lambda$ is given by [23]:

$$\lambda = \min(\lambda_{\min}, \max(\lambda_{\min}, \mathrm{trace}\ \hat{\boldsymbol{R}}_d(\omega_k))), \tag{9}$$

where $\hat{\boldsymbol{R}}_d(\omega_k)$ is an estimate of the target spatial covariance matrix, such as $\hat{\boldsymbol{R}}_d(\omega_k) = \hat{\boldsymbol{R}}_x(\omega_k) - \hat{\boldsymbol{R}}_v(\omega_k)$, and $\hat{\boldsymbol{R}}_v(\omega_k)$ is an estimate of the noise covariance matrix. The constants $\lambda_{\min}$ and $\lambda_{\max}$ are defined in Section 4. Diagonal loading constrains the norm of the weight vector and reduces the sensitivity of the beamformer to errors in the statistics [24]. This choice of $\lambda$ gives higher loading when the signal power is higher and smaller loading in noise-dominated segments.

## 3. PROPOSED ALGORITHM

Our goal is to improve the robust MWF introduced in Section 2.2 by exploiting frequency correlations in the target signal. To this end, we extend the narrowband model in (6) to form the *multiband* model, which incorporates frequency-shifted versions of the received signal. The frequency shifts are chosen so that the signal exhibits maximal self-correlation after shifting. Therefore, we focus on the fundamental frequency $\alpha_1$ and its integer multiples, known in acoustics as *harmonic* frequencies. The set $\mathcal{A}_1$ of modulation frequencies applied to the signal is defined as

$$\mathcal{A}_1 = \{\alpha_c : c\,\alpha_1, \ c = 0, \ldots, C-1\}, \tag{10}$$

where the number of modulations $C$ must be less than or equal to the number of harmonics in the signal, i.e., $C \leq P$. Next, we compute the STFT of each modulated signal and form a long vector $\boldsymbol{x}(\mathcal{A}_1, \omega_k) \in \mathbb{C}^{MC}$, by stacking the non-modulated noisy signal together with all modulated versions:

$$\boldsymbol{x}(\mathcal{A}_1, \omega_k)^T = [\boldsymbol{x}(\omega_k)^T \ \boldsymbol{x}(\omega_k - \alpha_1)^T \ \cdots \ \boldsymbol{x}(\omega_k - \alpha_{C-1})^T].$$

From here on, we write $\boldsymbol{x}(\mathcal{A}_1, \omega_k) = \boldsymbol{x}$ to represent multiband signals. The modulated reverberant signal vector $\boldsymbol{d}$ and the modulated noise vector $\boldsymbol{v}$ are constructed similarly, leading to:

$$\boldsymbol{x} = \boldsymbol{d} + \boldsymbol{v} \in \mathbb{C}^{MC}. \tag{11}$$

Now, notice that $\boldsymbol{d}$ can be represented by the matrix-vector multiplication $\boldsymbol{d} = \boldsymbol{C}\boldsymbol{s}$, where $\boldsymbol{s} = [\tilde{S}(\omega_k) \ \ldots \ \tilde{S}(\omega_k - \alpha_{C-1})]^T$ is the modulated signal at the reference microphone and $\boldsymbol{C} \in \mathbb{C}^{MC \times C}$ contains a frequency-shifted RTF padded with zeroes for each column. For example, for $C = 2$ we have

$$\boldsymbol{d} = \boldsymbol{C}\boldsymbol{s} = \begin{bmatrix} \boldsymbol{a}(\omega_k) & \boldsymbol{0}_{M(C-1)} \\ \boldsymbol{0}_{M(C-1)} & \boldsymbol{a}(\omega_k - \alpha_1) \end{bmatrix} \begin{bmatrix} \tilde{S}(\omega_k) \\ \tilde{S}(\omega_k - \alpha_1) \end{bmatrix}, \tag{12}$$

where $\boldsymbol{0}_A$ represents a zero vector of size $A$. Let us also define

$$\boldsymbol{S}_x(\mathcal{A}_1, \omega_k) = \boldsymbol{S}_x = \mathbb{E}[\boldsymbol{x}\boldsymbol{x}^H] \in \mathbb{C}^{MC \times MC} \tag{13}$$

as the spatial-spectral covariance matrix of the noisy signal. The spatial-spectral covariance matrices of the reverberant target and the noise are defined similarly and denoted by $\boldsymbol{S}_d$ and $\boldsymbol{S}_v$. Based on the multiband signal model, it is possible to extend the robust MWF

beamformer to optimally combine the noisy signals across different microphones and frequency shifts. The extended design is derived as the minimizer of the cost function below, which shares a similar form with (7):

$$J(\mathcal{A}_1, \omega_k) = \mathbb{E}[\|\tilde{S}(\omega_k) - \boldsymbol{w}^H \boldsymbol{x}\|_2^2] + \lambda \|\boldsymbol{w}\|_2^2. \qquad (14)$$

We use Wirtinger derivatives to obtain the gradient of (14) [25]. The solution is obtained by setting the gradient with respect to $\boldsymbol{w}^*$ to zero:

$$\boldsymbol{w}_{\text{cMWF}} = \boldsymbol{S}_\lambda^{-1} \boldsymbol{s}_{\boldsymbol{x}\tilde{s}} = \boldsymbol{S}_\lambda^{-1} \boldsymbol{s}_{\boldsymbol{d}\tilde{s}} = \boldsymbol{S}_\lambda^{-1} \boldsymbol{S}_{\boldsymbol{d}} \, \boldsymbol{e}_0, \qquad (15)$$

where we defined $\boldsymbol{S}_\lambda = \boldsymbol{S}_{\boldsymbol{x}} + \lambda \boldsymbol{I}$, $\boldsymbol{s}_{\boldsymbol{x}\tilde{s}} = \mathbb{E}[\boldsymbol{x} \tilde{S}^*(\omega_k)]$, $\boldsymbol{s}_{\boldsymbol{d}\tilde{s}} = \mathbb{E}[\boldsymbol{d} \tilde{S}^*(\omega_k)]$, $\boldsymbol{e}_0 = [1, 0, \ldots, 0]^T$ and the second equality follows from the assumption that the target and the noise are uncorrelated.

### 3.1. Estimating statistics

In practice, the harmonic frequencies of the signal and the spectral-spatial covariance matrices are unknown. The modulation set $\mathcal{A}_1$ is found by first estimating the fundamental frequency $\alpha_1$ using an algorithm based on non-linear least squares (NLS) [26]. The number of modulations $C$, which is related to the model order of the harmonic signal, is treated as a hyper-parameter and determined from the experiments. Section 3.2 provides additional details on how to handle the estimation of $\alpha_1$ on non-stationary data. The elements of $\boldsymbol{S}_{\boldsymbol{x}}$ are estimated using the ACP method detailed in Section 2.1 per each spectral frequency $\omega_k$, cyclic frequency $\alpha_c \in \mathcal{A}_1$ and microphone pair. $\boldsymbol{S}_{\boldsymbol{x}}$ is estimated from the noisy measurements, while $\boldsymbol{S}_{\boldsymbol{v}}$ is estimated from a noise-only segment. In most cases, the noise does not exhibit spectral correlation at the cyclic frequencies $\mathcal{A}_1$ associated with the target, i.e., only the null cyclic frequency $\alpha_0 = \{0\}$ is shared between the target and the noise. To enforce this assumption, we adjust the ACP estimate as $\hat{\boldsymbol{S}}_{\boldsymbol{v}} \leftarrow \text{blkdiag}(\hat{\boldsymbol{S}}_{\boldsymbol{v}})$, where $\text{blkdiag}(\cdot)$ extracts the block diagonal of a matrix, where each square block has size $M$. This modification retains only the spatial correlation of the noise while forcing its spectral correlation to 0, thereby reducing the number of unknowns in $\hat{\boldsymbol{S}}_{\boldsymbol{v}}$ from $M^2 C^2$ to $M^2 C$. The target $\hat{\boldsymbol{S}}_{\boldsymbol{d}}$ is then estimated using the generalized eigenvalue decomposition (GEVD) of $(\hat{\boldsymbol{S}}_{\boldsymbol{x}}, \hat{\boldsymbol{S}}_{\boldsymbol{v}})$, retaining only the $C$ eigenvectors associated with the $C$ largest eigenvalues, where $C$ is the maximum possible rank of $\boldsymbol{S}_{\boldsymbol{d}}$, since from the definition of $\boldsymbol{d}$ in (12), following the lines of [27, Lemma 1], we have:

$$\text{rank } \boldsymbol{S}_{\boldsymbol{d}} = \text{rank } \boldsymbol{C} \boldsymbol{S}_s \boldsymbol{C}^H \leq \min(\text{rank } \boldsymbol{C}, \text{rank } \boldsymbol{S}_s) \leq C.$$

The resulting blind cMWF is given by:

$$\hat{\boldsymbol{w}}_{\text{cMWF}} = \hat{\boldsymbol{S}}_\lambda^{-1} \hat{\boldsymbol{S}}_{\boldsymbol{d}}^{\text{gevd}} \boldsymbol{e}_0, \qquad (16)$$

where $\hat{\boldsymbol{S}}_\lambda = \hat{\boldsymbol{S}}_{\boldsymbol{x}} + \lambda \boldsymbol{I}$.

### 3.2. Recursive averaging

The estimation methods in Section 3.1 are valid as long as the fundamental frequency $\alpha_1$ remains fixed. However, the estimates $\hat{\alpha}_1(\ell)$ provided by the NLS algorithm varies from frame to frame. Direct use of $\hat{\alpha}_1(\ell)$ in computing the covariance matrices would require recomputing the modulated signals and their covariance matrices at every frame, yielding estimates with high variance. To mitigate this, define the relative temporal variation in the fundamental frequency at frame $\ell$ as:

$$\delta\alpha = |\hat{\alpha}_1(\ell) - \hat{\alpha}_1(\ell-1)| / (\hat{\alpha}_1(\ell-1) + \epsilon), \qquad (17)$$

where $\epsilon = 10^{-6}$ avoids divisions by 0 and $\hat{\alpha}_1(\ell) = 0$ if the $\ell$th frame is unvoiced. Next, introduce the smoothed fundamental frequency

estimate, $\bar{\alpha}_1(\ell)$, which is used to compute the time-dependent cyclic set $\mathcal{A}_1(\ell)$, the modulated signals, and their statistics:

$$\bar{\alpha}_1(\ell) = \begin{cases} \hat{\alpha}_1(\ell) & \text{if } D_0 \leq \delta\alpha < D_1, \\ \bar{\alpha}_1(\ell-1) & \text{otherwise}, \end{cases} \qquad (18)$$

where $D_0 < D_1$ are real-valued thresholds, and $\bar{\alpha}_1(0) = 0$. If $\delta\alpha < D_0$, the previous smoothed value is retained to avoid unnecessary re-modulations. If $\delta\alpha \geq D_1$, the previous value is retained because rapid variations would otherwise lead to poorly estimated statistics. When $\hat{\alpha}_1(\ell)$ changes moderately, $\bar{\alpha}_1(\ell)$ is updated accordingly.

At every frame, the current estimates of the spectral-spatial covariance matrices are updated with the new data. For example, $\hat{\boldsymbol{S}}_{\boldsymbol{x}}(\mathcal{A}_1(\ell), \ell)$ is updated as:

$$\hat{\boldsymbol{S}}_{\boldsymbol{x}}(\mathcal{A}_1(\ell), \ell) \leftarrow (1 - \beta)\hat{\boldsymbol{S}}_{\boldsymbol{x}}(\mathcal{A}_1(\ell-1), \ell-1) + \beta \, \boldsymbol{x}(\ell)\boldsymbol{x}(\ell)^H, \quad (19)$$

where the value of the constant $\beta$ is given in Section 4.2. Notice that the covariance matrices at different time frames may be functions of different modulation frequencies, thus the update is only approximately valid if the change in $\bar{\alpha}_1(\ell)$ is small and $C$ is the same. In Section 4, we see that if $\bar{\alpha}_1(\ell)$ changes slowly, as with the synthetic and instruments signals, the statistics are well estimated, and the cMWF performs well. For real speech, $\bar{\alpha}_1(\ell)$ and $C$ vary over time, complicating statistics estimation.

## 4. EXPERIMENTS

This section evaluates the proposed cMWF on simulated data, recordings from musical instruments, and speech signals. The blind beamformer in (16), which estimates the target covariance matrix through GEVD of $(\hat{\boldsymbol{S}}_{\boldsymbol{x}}, \hat{\boldsymbol{S}}_{\boldsymbol{v}})$, is compared against two unrealizable estimators that have access to ground-truth statistics. Results are given in terms SI-SDR improvements with respect to the unprocessed input at the first microphone [28]. The first oracle estimator, "cMWF+", uses the ACP estimate $\hat{\boldsymbol{S}}_{\boldsymbol{d}}$ of the ground-truth target instead of its GEVD estimate. It is given by:

$$\hat{\boldsymbol{w}}_{\text{cMWF}}^+ = (\hat{\boldsymbol{S}}_{\boldsymbol{d}} + \hat{\boldsymbol{S}}_{\boldsymbol{v}} + \lambda \boldsymbol{I})^{-1} \hat{\boldsymbol{S}}_{\boldsymbol{d}} \boldsymbol{e}_0. \qquad (20)$$

The second unrealizable estimator, termed "cMWF++", has access to the ACP estimate of the cross-statistics between the noisy and the ground truth target signals. It is given by:

$$\hat{\boldsymbol{w}}_{\text{cMWF}}^{++} = \hat{\boldsymbol{S}}_\lambda^{-1} \hat{\boldsymbol{s}}_{\boldsymbol{x}\tilde{s}}. \qquad (21)$$

To obtain similar variants of the narrowband MWFs, the spectral-spatial covariance matrices in (16), (20) and (21) are replaced by the corresponding SCMs. The amount of diagonal loading is calculated as in (9), with $\lambda_{\min} = 10^{-9}$ and $\lambda_{\max} = 10^{-4}$. When evaluating (9), we replace $\hat{\boldsymbol{R}}_{\boldsymbol{d}}$ with $\hat{\boldsymbol{S}}_{\boldsymbol{d}}$ for the cMWF variants. The cMWF variants are only used for frequency bins $\omega_k$ that lie close to the harmonic frequencies in $\mathcal{A}_1$, i.e., satisfying $|\omega_k - \alpha_c| < \varepsilon \Delta\omega$ for some $\alpha_c \in \mathcal{A}_1$, where $\Delta\omega$ is the spectral resolution and $\varepsilon$ is chosen as $\varepsilon = 1.5$. The remaining bins are processed with the narrowband MWF. We simulate a target point source, a directional interferer emitting white Gaussian noise (WGN) at $-10\,\text{dB}$ SNR measured at the reference microphone, and spatially uncorrelated WGN at $30\,\text{dB}$ SNR. The RIRs for the target and interferer are randomly selected from a set of 26 RIRs measured in a room with RT60 $= 0.61\,\text{s}$ and $8\,\text{cm}$ microphone spacing from the Bar-Ilan dataset [29]. Unless otherwise stated, we use $M = 2$ microphones, $C = 5$ frequency shifts, $K = 512$ points for the FFT, and the square-root Hann window with 75% overlap. The noise covariance matrix $\boldsymbol{S}_{\boldsymbol{v}}$ used in the GEVD is estimated from a separate, $2\,\text{s}$ long realization of the noise. Results are
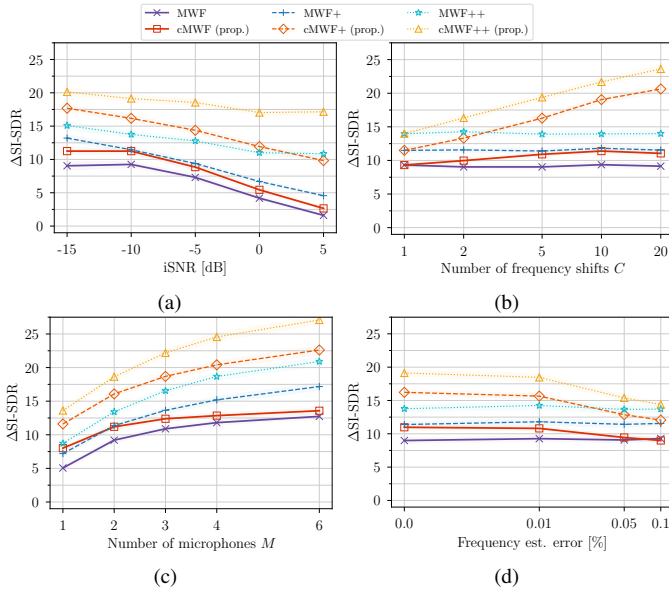
Fig. 1: Synthetic data. SI-SDR improvements over the noisy input for different beamformers. Each figure corresponds to a different varying parameter.

averaged over 50 Monte Carlo runs for the synthetic data experiments and over 10 runs for the real data experiments. We use different RIRs, noise, and target realizations in each run. The plot lines indicate the mean values, while shaded areas represent the 95% confidence intervals.

### 4.1. Synthetic data

First, we evaluate the accuracy of the beamformers on simulated data. The target signal is generated according to a simplified harmonic model [5, Eq. (7)], where the components at the different frequencies are perfectly correlated up to a multiplicative constant and a phase term. It is given by $Y(n) = B(n) \sum_{h=1}^{H} D_h \cos(\omega_0 nh + \phi_h)$, where $\{B(n)\}$ is a WSS process that describes the amplitude fluctuations over time, and the $D_h$ and the $\phi_h$ are random variables representing the relative amplitude and the phase of the sinusoids. The process $\{B(n)\}$ comprises independent Gaussian random variables distributed as $\mathcal{N}(0.5, 10)$ and lowpass filtered by a 4th order Butterworth filter with cutoff frequency $f_c = 5\,\mathrm{Hz}$. The $\phi_h$ are drawn from a uniform distribution $\mathcal{U}(-\pi, \pi)$, the $D_h$ are drawn from $\mathcal{U}(1, 10)$, and the frequency $\omega_0$ is drawn from $2\pi \cdot \mathcal{U}(60, 250)$. The number of harmonics $H$ is chosen as large as possible but obeys $\omega_0 H < f_s/2$. For the synthetic data experiments, the statistics are estimated using the entire signals, and $\alpha_1$ is assumed to be known. Each generated audio sample lasts 5 s. In Fig. 1a, we vary the input SNR due to the interferer (iSNR). Results indicate that the cyclic beamformers always improve performance over the conventional MWFs. Figure 1b shows how the SI-SDR increases with the number of frequency shifts $C$ in the cyclic models. Notice that, for $C = 20$ frequency shifts, "cMWF+" and "cMWF++" are approximately 10 dB SI-SDR better than "MWF+" and "MWF++", respectively. As mentioned earlier, if only $C = 1$ shift is considered, the cMWFs reduce to the corresponding MWFs. Next, we vary the number of microphones $M$ in Figure 1c. The performance of all beamformers improves when more microphones are available, as expected. Finally, Fig. 1d analyzes the sensitivity of the cMWF to a bias applied to the fundamental frequency used to compute the cyclic set $\mathcal{A}_1$ in (10). The perturbed fundamental frequency is given by $\dot{\alpha}_1 = \alpha_1(1 + \alpha_{\mathrm{err}}/100)$. The cyclic beamformers are only beneficial if the error in the fundamental frequency is less than 0.1%. Similarly, we found that the performance of the cMWFs degrades if
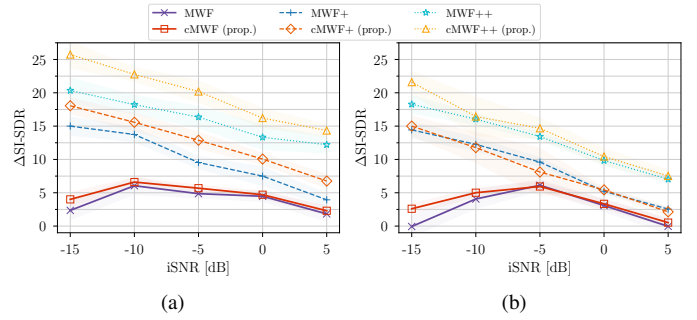


Fig. 2: Real data. SI-SDR improvements over the noisy input for the different beamformers. (a) shows results on the IOWA dataset and (b) on speech data.

the harmonics of the signal are not found at the exact integer multiples of the fundamental frequency (results not shown).

### 4.2. Real data

Next, we evaluate the recursive implementation of the algorithms described in Section 3.2 using music or speech recordings as targets. The first dataset comprises single-note brass instrument samples from [30]. When both vibrato and no-vibrato recordings are available, we select the latter ones. Only notes in the range C2 to C4 are considered, roughly corresponding to 65 Hz to 260 Hz. The second dataset consists of real speech recordings from the TIMIT database, uttered by either a male or a female speaker. In both experiments, we use $\beta = 0.05$ for recursive averaging of covariance matrices. In each Monte Carlo simulation, we randomly select 1 s of data. This value is chosen to be small because the single-note instrumental samples are of short duration. The constants that determine the update rate of $\bar{\alpha}_1(\ell)$ are set to $D_0 = 0.005$ and $D_1 = 0.2$. The fundamental frequency is estimated from the clean recordings. The cMWF variants are employed only when the fundamental frequency $\hat{\alpha}_1(\ell)$ has not changed significantly in the last frame to minimize the impact of poorly estimated covariance matrices; in other words, if $\delta\alpha \geq D_1$, we use the corresponding MWF variant for that time-frame. Improvements in SI-SDR are measured for different input SNRs and shown in Fig. 2. The cMWF variants consistently outperform the benchmark for instrument recordings (Fig. 2a), especially at lower iSNRs. For speech data (Fig. 2b), the blind cMWF has a better SI-SDR score for lower iSNRs, and it performs similarly to the benchmark for iSNR $-5$ dB or higher. PESQ [31] and STOI [32] results are omitted due to space limitations; they follow trends similar to SI-SDR. The non-blind variants of cMWF perform erratically on speech data. By analyzing the output spectrograms (not shown), we hypothesize that this is due to extremely large output values that sometimes occur when the fundamental frequency changes. In general, the reduced gains compared to Section 4.1 can be attributed to the limited accuracy in estimating $\alpha_1$ and the high variability of the fundamental frequency in speech. Additionally, whereas narrowband spatial covariance matrices change over time only by a scalar factor when sources are not moving, spectral covariance matrices vary with each note or phoneme, complicating their estimation.

### 5. CONCLUSION

This work proposed a cyclic multichannel Wiener filter for acoustic beamforming derived from a cyclostationary signal model. The beamformer exploits spectral correlation across harmonic frequencies in the target signal to enhance performance. While substantial SI-SDR gains are observed when the statistics, fundamental frequency, and harmonics are known exactly, as with synthetic data, the improvements are more modest on real recordings.

# REFERENCES

[1] S. M. Kay, *Fundamentals of Statistical Signal Processing: Estimation Theory*. Prentice-Hall, Inc., 1993.

[2] S. Gannot, E. Vincent, S. Markovich-Golan, and A. Ozerov, "A Consolidated Perspective on Multimicrophone Speech Enhancement and Source Separation," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2017.

[3] C. Li and R. C. Hendriks, "Alternating Least-Squares-Based Microphone Array Parameter Estimation for A Single-Source Reverberant and Noisy Acoustic Scenario," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2023.

[4] A. H. Moore, S. Hafezi, R. R. Vos, P. A. Naylor, and M. Brookes, "A Compact Noise Covariance Matrix Model for MVDR Beamforming," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2022.

[5] G. Bologni, R. Heusdens, and R. C. Hendriks, "Harmonics to the Rescue: Why Voiced Speech is Not a WSS Process," in *2024 18th International Workshop on Acoustic Signal Enhancement (IWAENC)*, 2024.

[6] W. A. Gardner, *Statistical Spectral Analysis: A Nonprobabilistic Theory*. Prentice-Hall, Inc., 1986.

[7] ——, *Cyclostationarity in Communications and Signal Processing*. IEEE Press, 1994.

[8] W. A. Gardner, A. Napolitano, and L. Paura, "Cyclostationarity: Half a century of research," *Signal Processing*, 2006.

[9] B. Feher, "Short Overview of Cyclostationary Signal Processing." Technische Univ. Delft (Netherlands). Dept. of Applied Physics., Tech. Rep., 1995.

[10] W. Gardner, "Cyclic Wiener filtering: Theory and method," *IEEE Trans. Commun.*, Jan./1993.

[11] A. Napolitano, "A - Nonstationary signal analysis," in *Cyclostationary Processes and Time Series*. Academic Press, 2020.

[12] B. Moore, *Hearing*. Academic Press, 1995.

[13] H. Zhang, A. Abdi, and A. Haimovich, "Reduced-Rank Multi-Antenna Cyclic Wiener Filtering for Interference Cancellation," in *MILCOM 2006 - 2006 IEEE Military Communications Conference*, 2006.

[14] G. Bologni, "Cyclic multichannel wiener filter," [Online]. Available at https://github.com/Screeen/cyclicMWF.

[15] W. Gardner, "Measurement of spectral correlation," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 1986.

[16] J. Antoni, "Cyclic spectral analysis of rolling-element bearing signals: Facts and fictions," *Journal of Sound and Vibration*, 2007.

[17] R. Roberts, W. Brown, and H. Loomis, "Computationally efficient algorithms for cyclic spectral analysis," *IEEE Signal Processing Magazine*, 1991.

[18] P. Borghesani and J. Antoni, "A faster algorithm for the calculation of the fast spectral correlation," *Mechanical Systems and Signal Processing*, 2018.

[19] J. K. Alsalaet, "Fast Averaged Cyclic Periodogram method to compute spectral correlation and coherence," *ISA Transactions*, 2022.

[20] S. Braun, A. Kuklasiński, O. Schwartz, O. Thiergart, E. A. P. Habets, S. Gannot, S. Doclo, and J. Jensen, "Evaluation and Comparison of Late Reverberation Power Spectral Density Estimators," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2018.

[21] S. Doclo, A. Spriet, J. Wouters, and M. Moonen, "Speech Distortion Weighted Multichannel Wiener Filtering Techniques for Noise Reduction," in *Speech Enhancement*. Springer, 2005.

[22] J. Li, P. Stoica, and Z. Wang, "On robust Capon beamforming and diagonal loading," *IEEE Transactions on Signal Processing*, 2003.

[23] Q. Wu and K. M. Wong, "Blind adaptive beamforming for cyclostationary signals," *IEEE Transactions on Signal Processing*, Nov./1996.

[24] H. Cox, R. Zeskind, and M. Owen, "Robust adaptive beamforming," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 1987.

[25] D. H. Brandwood, "A complex gradient operator and its application in adaptive array theory," *IEE Proceedings F (Communications, Radar and Signal Processing)*, 1983.

[26] J. K. Nielsen, T. L. Jensen, J. R. Jensen, M. G. Christensen, and S. H. Jensen, "Fast fundamental frequency estimation: Making a statistically efficient estimator computationally efficient," *Signal Processing*, 2017.

[27] G. Bologni, R. C. Hendriks, and R. Heusdens, "Wideband Relative Transfer Function (RTF) Estimation Exploiting Frequency Correlations," *IEEE Transactions on Audio, Speech and Language Processing*, 2025.

[28] J. L. Roux, S. Wisdom, H. Erdogan, and J. R. Hershey, "SDR – Half-baked or Well Done?" in *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019.

[29] E. Hadad, F. Heese, P. Vary, and S. Gannot, "Multichannel audio database in various acoustic environments," in *2014 14th International Workshop on Acoustic Signal Enhancement (IWAENC)*, 2014.

[30] University of Iowa Electronic Music Studios, "Musical instrument samples," 2011, [Online]. Available at https://theremin.music.uiowa.edu/MIS.html.

[31] A. Rix, J. Beerends, M. Hollier, and A. Hekstra, "Perceptual evaluation of speech quality (PESQ)-a new method for speech quality assessment of telephone networks and codecs," in *2001 IEEE International Conference on Acoustics, Speech, and Signal Processing. Proceedings (Cat. No.01CH37221)*, 2001.

[32] C. H. Taal, R. C. Hendriks, R. Heusdens, and J. Jensen, "An Algorithm for Intelligibility Prediction of Time–Frequency Weighted Noisy Speech," *IEEE Transactions on Audio, Speech, and Language Processing*, 2011.