# LOW-RANKNESS OF COMPLEX-VALUED SPECTROGRAM AND ITS APPLICATION TO PHASE-AWARE AUDIO PROCESSING

*Yoshiki Masuyama, Kohei Yatabe and Yasuhiro Oikawa*

Department of Intermedia Art and Science, Waseda University, Tokyo, Japan

## ABSTRACT

Low-rankness of amplitude spectrograms has been effectively utilized in audio signal processing methods including non-negative matrix factorization. However, such methods have a fundamental limitation owing to their amplitude-only treatment where the phase of the observed signal is utilized for resynthesizing the estimated signal. In order to address this limitation, we directly treat a complex-valued spectrogram and show a complex-valued spectrogram of a sum of sinusoids can be approximately low-rank by modifying its phase. For evaluating the applicability of the proposed low-rank representation, we further propose a convex prior emphasizing harmonic signals, and it is applied to audio denoising.

***Index Terms***— Instantaneous frequency, phase derivative, data-driven approach, convex-optimization, audio denoising.

## 1. INTRODUCTION

In audio signal processing, low-rank representation of amplitude spectrograms has been utilized extensively [1–7]. For instance, the non-negative matrix factorization (NMF) has been developed in many applications including source separation [8], audio inpainting [9], and music transcription [10]. While those methods have successfully applied to various problems, they have a limitation owing to their amplitude-only treatment.

Recent studies have shown the importance of phase [11–13], and some phase-aware extensions of NMF have been proposed [14, 15]. In addition to amplitude spectrograms, the complex NMF (CNMF) treats phase at each time-frequency bin as the independent variables to be optimized [14]. Meanwhile, the time-domain spectrogram factorization (TSF) implements NMF-like signal decomposition in the time domain for implicitly considering phase based on the consistency of a spectrogram [15]. Although these methods include phase in their models, low-rankness is only imposed on amplitude spectrograms, and the explicit structure of the phase was not considered.

A very recent study revealed the relation between the rank and phase of a complex-valued spectrogram [16]. The phase of a sum of sinusoids has a distinctive structure which has been widely utilized in audio signal processing [17–22], and the rank of its complex-valued spectrogram depends on the number of the sinusoids [16]. This theoretical result indicates that the rank of the complex-valued spectrogram increases as the number of sinusoids increases, while its amplitude can stay low-rank. Therefore, imposing low-rankness only on amplitude spectrograms as in CNMF and TSF seems to be justified, and low-rank treatment of a complex-valued spectrogram sounds inappropriate. However, we found the complex-valued spectrogram can be well approximated by a low-rank matrix applying a specific modification of phase.

In this paper, we propose a low-rank representation of a complex-valued spectrogram by applying the instantaneous phase correction introduced in [19]. This phase modification is based on the phase model of harmonic signals, and the phase corrected complex-valued spectrogram of harmonic signals becomes as low-rank as its amplitude spectrogram with some assumptions. As an example of the applications of the proposed low-rank representation, we further propose a convex prior which emphasizes sinusoidal components through the low-rankness, and its effectiveness is demonstrated by an audio denoising experiment.

## 2. PREVIOUS WORK

Let the short-time Fourier transform (STFT) of a signal $\mathbf{x} = [x[0], \ldots, x[L-1]]^\top \in \mathbb{R}^L$ with a window $\mathbf{w} \in \mathbb{R}^L$ be

$$\mathcal{G}^{\mathbf{w}}(\mathbf{x})[\xi, \tau] = \sum_{l=0}^{L-1} x[l + a\tau]w[l]e^{-2\pi j\xi bl/L}, \quad (1)$$

where $j = \sqrt{-1}$, $a$ and $b$ are the time and frequency shifting steps, $\tau = 0, 1, \ldots, T-1$ and $\xi = 0, 1, \ldots, K-1$ denote the time and frequency indices, and index overflow is treated by zero-padding. This STFT can be represented by a matrix form:

$$\mathcal{G}^{\mathbf{w}}(\mathbf{x}) = \mathbf{F}\text{diag}(\mathbf{w})\mathbf{X}, \quad (2)$$

$$\mathbf{X} = [\mathbf{x}_0, \mathbf{x}_1, \ldots, \mathbf{x}_{T-1}], \quad (3)$$

where $\mathbf{X}$ is the horizontal concatenation of $T$ patches of the signal $\mathbf{x}_\tau = [x[a\tau], x[a\tau+1], \ldots, x[a\tau+L-1]]^\top$, $\mathbf{z}^\top$ is the transpose of $\mathbf{z}$, $\mathbf{F}$ is the discrete Fourier transform matrix, and $\text{diag}(\mathbf{w})$ denotes the diagonal matrix whose diagonal elements are given by $\mathbf{w}$.

Let a sinusoid be written as

$$s_0[l] = A_0 e^{2\pi j f_0 l/L + \phi_0}, \quad (4)$$

where $A_0 \in \mathbb{R}_+$, $f_0 \in [0, L/2)$ and $\phi_0 \in [0, 2\pi)$ are the amplitude, frequency, and initial phase of the sinusoid, respectively. Considering patches of this sinusoid $\mathbf{s}_0^0, \mathbf{s}_1^0, \ldots, \mathbf{s}_{T-1}^0$ given by

$$\mathbf{s}_\tau^0 = [s_0[a\tau], s_0[a\tau+1], \ldots, s_0[a\tau+L-1]]^\top, \quad (5)$$

every patch has the following relation:

$$\mathbf{s}_\tau^0 = e^{2\pi j f a/L}\mathbf{s}_{\tau-1}^0 = e^{2\pi j f a\tau/L}\mathbf{s}_0^0. \quad (6)$$

Considering the matrix $\mathbf{S}_0$ whose columns are given by $\mathbf{s}_\tau^0$ as Eq. (3), $\text{rank}(\mathbf{S}_0) = 1$, where $\text{rank}(\mathbf{S}_0)$ is the rank of the matrix $\mathbf{S}_0$. This is because all columns of $\mathbf{S}_0$ are given by a complex-valued scalar multiplication of $\mathbf{s}_0$. Since $\text{rank}(\mathcal{G}^{\mathbf{w}}(\mathbf{x}))$ coincides with $\text{rank}(\mathbf{X})$ when $\text{diag}(\mathbf{w})$ is full rank (i.e., $w[l] \neq 0 \ \forall l$) by the unitarity of $\mathbf{F}$, the rank of the complex-valued spectrogram of a sinusoid is 1 as described in [16].

The previous work also showed that the rank of a complex-valued spectrogram of a sum of sinusoids depends on the number of sinusoids [16], while its amplitude can be well approximated by a

rank-1 matrix regardless of the number of sinusoids. Consider a sum of $H$ sinusoids:

$$s[l] = \sum_{h=0}^{H-1} s_h[l] = \sum_{h=0}^{H-1} A_h e^{2\pi j f_h l/L + \phi_h}, \quad (7)$$

where $A_h \in \mathbb{R}_+$, $f_h \in [0, L/2)$ ($f_p \neq f_q$ when $p \neq q$), and $\phi_h \in [0, 2\pi)$ are the amplitude, frequency, and initial phase of $h$th sinusoid, respectively. The matrix which contains its patches is defined as $\mathbf{S} = \sum_{h=0}^{H-1} \mathbf{S}_h$, and $\mathrm{rank}(\mathcal{G}^{\mathbf{w}}(\mathbf{s}))$ increases as the number of sinusoids $H$ increases[1].

Low-rank representation of an amplitude spectrogram has been well-accepted in audio signal processing since an amplitude spectrogram of a sum of sinusoids is low-rank regardless of the number of sinusoids. In contrast, the low-rankness of a complex-valued spectrogram has not been considered owing to the above nature.

## 3. PROPOSED LOW-RANK REPRESENTATION OF COMPLEX-VALUED SPECTROGRAMS

In this section, we show that a complex-valued spectrogram of a sum of sinusoids becomes low-rank when the instantaneous phase correction [19] is applied. The characteristic of the complex-valued spectrogram with this phase correction is reviewed first, and then its low-rankness is validated with some numerical examples.

### 3.1. Instantaneous phase corrected STFT (iPC-STFT) [19]

As in the previous section, the rank of a complex-valued spectrogram of a sum of sinusoids depends on the number of sinusoids, while its amplitude is well represented by a rank-1 matrix. This is because the evolution of the phase is different for each sinusoid as illustrated in Fig. 1(a), where a sum of two sinusoids is shown in the time-frequency domain. Since the phase evolves along time, the real and imaginary parts of the complex-valued spectrogram periodically fluctuate. Their periods are different in each sinusoid, which results in the increase of the rank. If this phase evolution is eliminated, the rank of the complex-valued spectrogram can be reduced.

The phase evolution is closely related to the instantaneous frequency at each time-frequency bin. When we assume each sinusoid in Eq. (7) is sufficiently separated (i.e., the interference from other sinusoids can be ignored in the region dominated by a sinusoid), the complex-valued spectrogram of the sum of sinusoids has the following relation by utilizing the instantaneous frequency [21]:

$$\mathcal{G}^{\mathbf{w}}(\mathbf{s})[\xi, \tau+1] = e^{2\pi j v[\xi,\tau]a/L} \mathcal{G}^{\mathbf{w}}(\mathbf{s})[\xi, \tau], \quad (8)$$

where $v[\xi, \tau]$ is different for each sub-band but same in all time-frames. That is, the phase of the complex-valued spectrogram evolves as a constant multiple of the instantaneous frequency $v[\xi, \tau]$.

To cancel this phase evolution, the instantaneous phase corrected STFT (iPC-STFT) was proposed as follows: [19]

$$\mathcal{G}_{\mathrm{iPC}}^{\mathbf{w}}(\mathbf{x}) = \mathbf{E} \odot \mathcal{G}^{\mathbf{w}}(\mathbf{x}), \quad (9)$$

where $\odot$ is the Hadamard product, $\mathbf{E}$ is the instantaneous phase correction matrix whose element is defined by

$$E[\xi, \tau] = \prod_{\eta=0}^{\tau-1} e^{-2\pi j v[\xi,\eta]a/L}, \quad (10)$$

---

[1] As described in [16], the complex-valued spectrogram of a sum of $H$ sinusoids becomes a rank-$H$ matrix when the frequencies of sinusoids are on the discrete Fourier grid. Note that [16] considers STFT with maximal redundancy and periodic extension of the signal and window. An exact characterization for more general cases was not presented.
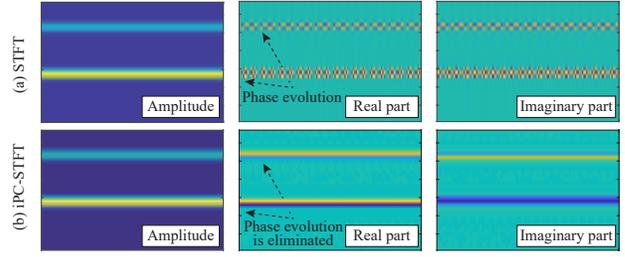


**Fig. 1**. Comparison of spectrograms of two sinusoids calculated by (a) the usual STFT and (b) iPC-STFT.

and $E[\xi, 0] = 1$ for all $\xi$. This matrix cancels the phase evolution in Eq. (8), which results in the low-rank complex-valued spectrogram as shown in Fig. 1(b). Both real and imaginary parts of the complex-valued spectrogram are constant at each sub-band, and therefore its rank is reduced. Note that this correction of phase can be easily inverted by multiplying the complex conjugate of $\mathbf{E}$.

In reality, the instantaneous frequency $v[\xi, \tau]$ is not known and must be estimated from the observed signal. One simple method is to directly calculate the time-differential of phase:

$$v[\xi, \tau] = b\xi - \mathrm{Im}\left[\frac{\mathcal{G}^{\mathbf{w}'}(\mathbf{x})[\xi, \tau]}{\mathcal{G}^{\mathbf{w}}(\mathbf{x})[\xi, \tau]}\right], \quad (11)$$

where $\mathbf{w}'$ is time-derivative of the window $\mathbf{w}$ [23–25], and $\mathrm{Im}[z]$ is the imaginary part of $z$. Once the instantaneous frequency is estimated, the instantaneous phase correction is uniquely defined as an invertible linear transform in the proposed denoising scheme introduced in Section 4.1. Note that iPC-STFT can be applied to arbitrary signal which may not consist of pure sinusoids.

### 3.2. Low-rankness of iPC-STFT spectrogram

The previous work [19] only considered time-directional smoothness of iPC-STFT, and no further characteristics have been shown. In this paper, we show the low-rankness of the complex-valued spectrogram calculated by iPC-STFT. Such low-rankness in the complex domain should be important for extending ordinary studies of low-rank audio modeling.

From Eqs. (8)–(10), with some conditions mentioned in the previous subsection, the following neighborhood relation can be obtained for iPC-STFT of a sum of sinusoids:

$$\mathcal{G}_{\mathrm{iPC}}^{\mathbf{w}}(\mathbf{s})[\xi, \tau+1] = \prod_{\eta=0}^{\tau} e^{-2\pi j v[\xi,\eta]a/L} \mathcal{G}^{\mathbf{w}}(\mathbf{s})[\xi, \tau+1],$$

$$= \prod_{\eta=0}^{\tau-1} e^{-2\pi j v[\xi,\eta]a/L} \mathcal{G}^{\mathbf{w}}(\mathbf{s})[\xi, \tau],$$

$$= \mathcal{G}^{\mathbf{w}}(\mathbf{s})[\xi, 0] \quad (= \mathcal{G}_{\mathrm{iPC}}^{\mathbf{w}}(\mathbf{s})[\xi, 0]). \quad (12)$$

That is, all columns of $\mathcal{G}_{\mathrm{iPC}}^{\mathbf{w}}(\mathbf{s})$ is equivalent to the first column of $\mathcal{G}_{\mathrm{iPC}}^{\mathbf{w}}(\mathbf{s})$ when the instantaneous phase correction eliminates the phase evolution completely (i.e., when the instantaneous frequency is accurately estimated). Hence, $\mathrm{rank}(\mathcal{G}_{\mathrm{iPC}}^{\mathbf{w}}(\mathbf{s})) = 1$ regardless of the number of sinusoids.

The above relation of iPC-STFT can be generalized to a signal beyond a sum of sinusoids. The relation in Eq. (8) can be understood as the well-accepted sinusoidal model [18]:

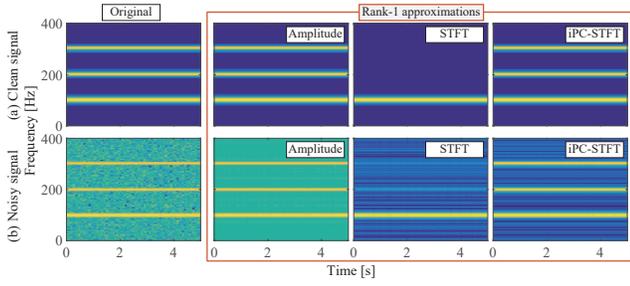$$\phi[\xi, \tau+1] = \phi[\xi, \tau] + 2\pi v[\xi, \tau]a/L, \quad (13)$$

**Fig. 2**. Amplitude of the rank-1 approximated spectrograms of a clean sum of three sinusoids and degraded one.

**Table 1**. SNR of rank-1 approximation of spectrograms in Fig. 2.

| | shift size | \multicolumn{4}{c}{Input SNR [dB]} | | | |
| --- | --- | --- | --- | --- | --- |
| | shift size | 0 | 10 | 20 | Clean |
| Amplitude | 1/2 | 1.3 | 11.3 | 21.4 | 64.4 |
| | 1/4 | 1.3 | 11.4 | 21.4 | 64.3 |
| | 1/8 | 1.3 | 11.4 | 21.4 | 62.9 |
| STFT | 1/2 | 2.2 | 2.3 | 2.3 | 2.3 |
| | 1/4 | 2.2 | 2.3 | 2.3 | 2.3 |
| | 1/8 | 2.3 | 2.3 | 2.3 | 2.3 |
| iPC-STFT | 1/2 | 18.8 | 28.9 | 38.7 | 52.3 |
| | 1/4 | 21.8 | 31.6 | 41.5 | 55.4 |
| | 1/8 | 24.5 | 34.3 | 44.2 | 55.3 |

where $\phi$ is a phase spectrogram. This equation indicates that $\phi[\xi,\tau] = \phi[\xi,0] + \sum_{\eta=0}^{\tau-1} 2\pi v[\xi,\eta]a/L$, and therefore

$$
\begin{aligned}
\mathcal{G}_{\mathrm{iPC}}^{\mathbf{w}}(\mathbf{x})[\xi,\tau] &= \prod_{\eta=0}^{\tau-1} e^{-2\pi j v[\xi,\eta]a/L} \mathcal{G}^{\mathbf{w}}(\mathbf{x})[\xi,\tau], \\
&= \prod_{\eta=0}^{\tau-1} e^{-2\pi j v[\xi,\eta]a/L} |\mathcal{G}^{\mathbf{w}}(\mathbf{x})[\xi,\tau]| e^{2\pi j \phi[\xi,\tau]a/L}, \\
&= |\mathcal{G}^{\mathbf{w}}(\mathbf{x})[\xi,\tau]| e^{2\pi j \phi[\xi,0]a/L}, \qquad (14)
\end{aligned}
$$

i.e., the instantaneous phase correction converts a complex-valued spectrogram into a constant multiple of its amplitude. Hence, $\mathrm{rank}(\mathcal{G}_{\mathrm{iPC}}^{\mathbf{w}}(\mathbf{x})) = \mathrm{rank}(|\mathcal{G}^{\mathbf{w}}(\mathbf{x})|)$ whenever the instantaneous frequency at each time-frequency bin $v[\xi,\tau]$ is estimated exactly. When $v[\xi,\tau]$ can be estimated only approximately, the relation also becomes approximation: $\mathrm{rank}(\mathcal{G}_{\mathrm{iPC}}^{\mathbf{w}}(\mathbf{x})) \approx \mathrm{rank}(|\mathcal{G}^{\mathbf{w}}(\mathbf{x})|)$.

We stress that the low-rankness of a complex-valued spectrogram is completely different from that of the amplitude. For instance, amplitude spectrogram of the Gaussian noise realized in the time domain can be well approximated by rank-1 because its energy is almost the same for all time-frequency bins. In contrast, its complex-valued spectrogram is not low-rank because its phase dose not obey Eq. (13). Such difference between the low-rankness of amplitude and complex-valued spectrograms is experimentally illustrated in the next subsection.

### 3.3. Numerical examples

For illustrating the property of iPC-STFT and its low-rank representation, some simple examples are shown here. Firstly, a sum of three sinusoids, $s[l] = \sum_{h=0}^{H-1} A_h \sin(2\pi f_h l/L)$, is considered, where $H = 3$, $A_h = 10 - h$, $f_h = (h+1)f_0$, $f_0 = 100$ Hz, and the sampling frequency was 16000 Hz. STFT was calculated with the Hann
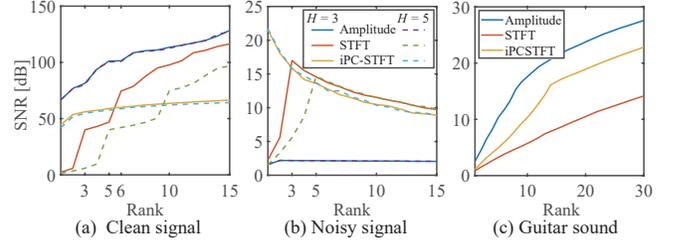


**Fig. 3**. SNR of rank-$k$ approximation of spectrograms of (a) a sum of pure sinusoids, (b) (a) with the additive Gaussian noise, and (c) the guitar sound without noise.

window (4096 samples), and the instantaneous frequency was estimated by Eq. (11). Here the rank-$k$ approximations were calculated by the truncated singular value decomposition.

Rank-1 approximations of the three sinusoids are illustrated in Fig. 2. The leftmost column is the original signals, and second to fourth columns represent the rank-1 approximations of the amplitude, complex (STFT), and complex (iPC-STFT) spectrograms, respectively. From Fig. 2(a), the usual complex-valued (STFT) spectrogram can only represent one sinusoid by the rank-1 approximation as described in Section 2. In contrast, both amplitude and iPC-STFT spectrograms can simultaneously represent all three sinusoids by rank-1 approximation. This result confirms the relation $\mathrm{rank}(|\mathcal{G}^{\mathbf{w}}(\mathbf{x})|) \approx \mathrm{rank}(\mathcal{G}_{\mathrm{iPC}}^{\mathbf{w}}(\mathbf{x}))$ described in the previous subsection. Fig. 2(b) represents a noisy version of Fig. 2(a), where the complex Gaussian noise was added in the time-frequency domain. In this noisy case, the rank-1 approximation of the amplitude spectrogram stayed noisy because the amplitude spectrogram of stationary noise can be well approximated by a low-rank matrix. On the other hand, rank-1 approximation of complex-valued (STFT and iPC-STFT) spectrograms can remove the Gaussian noise to some extent. This is because iPC-STFT takes the phase structure of sinusoidal components into account, while the amplitude specrogram ignores the phase information. That is, the proposed low-rank representation of complex-valued spectrograms can distinguish sinusoidal components from noise. Some quantitative data of those rank-1 approximations are shown in Table 1, where the signal-to-noise ratio (SNR) of the original signal and the shift size of window are varied. It confirms that only iPC-STFT can improve SNR of Fig. 2(b) by the rank-1 approximation[2].

Next, SNR of rank-$k$ approximations of sums of three or five sinusoids are shown in Fig. 3(a) and 3(b). The shift size of the window was 1/4, and the input SNR was 10 dB for Fig. 3(b). The rank-$k$ approximation of the amplitude spectrogram resulted in lowest SNR because it does not distinguish sinusoidal components from the stationary noise. Since the rank of the usual complex-valued (STFT) spectrogram is decided by the number of sinusoidal components [16], SNR of STFT was low for the rank-1 approximation ($k = 1$) but significantly improved when the rank $k$ coincides with the number of sinusoids $H$. For the noisy situation in Fig. 3(b), the highest SNR of STFT was achieved when at $k = H$. In contrast, iPC-STFT obtained a low-rank complex-valued spectrograms which can be well approximated by a rank-1 matrix. As a result, the highest SNR in Fig. 3(b) was achieved by the rank-1 approximation of the complex-valued iPC-STFT spectrogram regardless of the number of sinusoids $H$. Finally, rank-$k$ approximation of a real guitar sound is illustrated in Fig. 3(c), where the guitar sample was ob-

---

[2] For calculating SNR of the rank-1 approximation of amplitude spectrograms in the complex domain, the observed noisy phase was utilized. This approach was often utilized in NMF for resynthesizing the estimated signal.

tained from IDMT-SMT-GUITAR database[3] [26], and the sampling frequency was 44100 Hz. The complex-valued spectrogram calculated by iPC-STFT can represent the real guitar sound better than that of the usual STFT. One reason of lower SNR of iPC-STFT comparing to the amplitude-only approximation should be the error of the instantaneous frequency estimation which can be improved by considering an estimation method more sophisticated than Eq. (11).

## 4. APPLICATION

### 4.1. Proposed low-rank representation as a signal prior

For applying the proposed low-rank representation to audio signal processing, we propose a novel prior, named instantaneous phase corrected low-rankness (iPCLR), defined by

$$\mathcal{P}_{\mathrm{iPCLR}}(\mathbf{x}) = \|\mathcal{G}_{\mathrm{iPC}}^{\mathbf{w}}(\mathbf{x})\|_* , \qquad (15)$$

where $\| \cdot \|_*$ denotes the nuclear norm which is the convex envelope of the rank function $\mathrm{rank}(\cdot)$. Thanks to the property of iPC-STFT illustrated in the previous section, this prior emphasizes sinusoidal components and suppresses non-sinusoidal components such as random noise. Note that the instantaneous frequency is pre-computed and fixed to make iPC-STFT into a linear operator (independent of $\mathbf{x}$), and thus this prior is convex and can be easily incorporated with other priors.

As an application of iPCLR, the following convex optimization problem is considered for audio denoising,

$$\mathbf{x}^\star = \arg\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{x} - \mathbf{d}\|_2^2 + \lambda\,\mathcal{P}_{\mathrm{iPCLR}}(\mathbf{x}), \qquad (16)$$

where $\mathbf{d}$ is the noisy input signal, and $\lambda > 0$ is a regularization parameter. Here, the instantaneous phase correction matrix $\mathbf{E}$ for iPC-STFT is fixed to that calculated from the noisy input signal, and thus iPC-STFT is a fixed linear operator. This simple denoiser is called *proximity operator* which is a building block of the optimization-based signal processing for solving a variety of problems [27]. That is, if Eq. (16) is effective for denoising, then the proposed prior $\mathcal{P}_{\mathrm{iPCLR}}$ should be effective in other applications as well.

### 4.2. Experimental result of audio denoising

The proposed method in Eq. (16) was applied to audio denoising of the three melodies played by different musical instruments from songKitamura dataset [28], where the Gaussian noise was added in the time domain. The sampling frequency was 44100 Hz, and STFT was calculated by the canonical tight window of the Hann window of 4096 samples with 1024 sample shifting.

The proposed method was compared with other low-rank models: Euclidean NMF (EUC-NMF) for amplitude spectrograms, Itakura–Saito NMF (IS-NMF) [1] for power spectrograms, and CNMF [14]. TSF [15] was also compared with a slight modification, the perfect reconstruction constraint $\mathbf{d} = \sum_h \mathbf{x}_h$ was relaxed to a penalty $\beta/2\|\mathbf{d} - \sum_h \mathbf{x}_h\|_2^2$, because the original formulation of TSF is not suitable for a denoising application. The number of bases was set to 30 for the conventional methods (note that the degree of low-rankness is decided by $\lambda$ in the proposed method). The other parameters of CNMF and TSF were set to the default value in the original papers [14, 15]. The number of iterations was set to 100 for all methods, where ADMM [29] was adopted for solving Eq. (16), and 10 initial values were randomly chosen. The
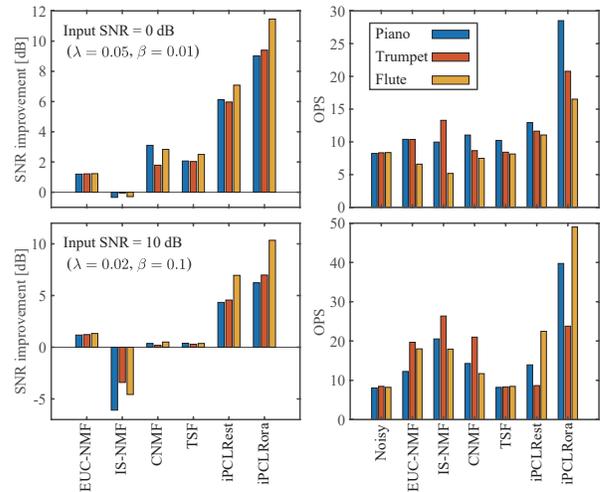
---

[3]http://www.idmt.fraunhofer.de/en/business_units/smt/guitar.html



**Fig. 4**. SNR and OPS of the denoising results. They were the average scores of 10 trials whose initial values were randomly selected.

proposed iPCLR was implemented in two ways for investigating the effect of estimation error on the instantaneous frequency; iPCLRest estimated the instantaneous frequency from the noisy signal, while iPCLRora estimated it from the original clean signal. In both cases, the estimation was done by Eq. (11). Note that the proposed method and TSF are phase-aware in the sense that the variable is treated in the time domain which ensures consistency [15] of the spectrogram.

Fig. 4 shows the average score calculated from ten random initial values, where the performances were evaluated by SNR and overall-perceptual-score (OPS) which is a perceptual measure available in the PEASS toolbox [30]. The low-rank models of amplitude spectrograms (EUC- and IS-NMF) obtained limited results as similar to those in the table and figures in the previous section. This is because those models treat the stationary noise as a part of the low-rank components and do not attempt to reduce it. For noisier situation (top row), CNMF and TSF obtained higher SNR comparing to EUC- and IS-NMF by considering phase. On the other hand, their SNR improvement were lower than that of EUC-NMF when the input SNR $= 10$ dB (bottom row), which should be because they do not consider the explicit structure of phase that causes instability. In contrast to the conventional methods, the proposed method in Eq. (16) achieved better scores by taking advantage of considering the structure of the phase given by Eq. (13), even when the instantaneous frequency was estimated from the noisy observations (iPCLRest). Thanks to the accurate estimation of the instantaneous frequency, iPCLRora resulted in the highest SNR and OPS. However, we stress that iPCLRest also worked well in terms of SNR because the instantaneous frequency around the spectral peaks can be accurately estimated. The error of the instantaneous frequency at the time-frequency bin with small amplitude does not significantly affect to the proposed method.

## 5. CONCLUSION

In this paper, we showed that the rank of a complex-valued spectrogram can be as low as its amplitude by applying the instantaneous phase correction under mild assumptions. Based on this finding, a low-rank model called iPCLR was proposed for audio signal processing, and its potentiality was illustrated through audio denoising. Seeking further applications of iPCLR remains as future works.

*This paper has been accepted to the 44th International Conference on Acoustics, Speech, and Signal Processing (ICASSP 2019).*

## 6. REFERENCES

[1] C. Févotte, N. Bertin, and J.-L. Durrieu, "Nonnegative matrix factorization with the itakura-saito divergence: With application to music analysis," *Neural Comput.*, vol. 21, no. 3, pp. 793–830, Mar. 2009.

[2] A. Ozerov and C. Févotte, "Multichannel nonnegative matrix factorization in convolutive mixtures for audio source separation," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 18, no. 3, pp. 550–563, Mar. 2010.

[3] K. W. Wilson, B. Raj, and P. Smaragdis, "Regularized nonnegative matrix factorization with temporal dependencies for speech denoising," in *INTERSPEECH*, Sept. 2008, pp. 411–414.

[4] P. S. Huang, S. D. Chen, P. Smaragdis, and M. Hasegawa-Johnson, "Singing-voice separation from monaural recordings using robust principal component analysis," in *IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, Mar. 2012, pp. 57–60.

[5] T. Komatsu, Y. Senda, and R. Kondo, "Acoustic event detection based on non-negative matrix factorization with mixtures of local dictionaries and activation aggregation," in *IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, Mar. 2016, pp. 2259–2263.

[6] D. Kitamura, N. Ono, H. Sawada, H. Kameoka, and H. Saruwatari, "Determined blind source separation unifying independent vector analysis and nonnegative matrix factorization," *IEEE/ACM Trans. Audio, Speech Lang. Process*, vol. 24, no. 9, pp. 1622–1637, Sept. 2016.

[7] K. Yatabe and D. Kitamura, "Determined blind source separation via proximal splitting algorithm," in *IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, Apr. 2018, pp. 776–780.

[8] T. Virtanen, "Monaural sound source separation by nonnegative matrix factorization with temporal continuity and sparseness criteria," *IEEE Trans. Audio, Speech, and Lang. Process.*, vol. 15, no. 3, pp. 1066–1074, Mar. 2007.

[9] P. Smaragdis, B. Raj, and M. Shashanka, "Missing data imputation for spectral audio signals," in *IEEE Int. Workshop Mach. Learn. Signal Process.*, Sept. 2009, pp. 1–6.

[10] M. D. Plumbley, S. A. Abdallah, J. P. Bello, M. E. Davies, G. Monti, and M. B. Sandler, "Automatic music transcription and audio source separation," *Cybern. Syst.*, vol. 33, no. 6, pp. 603–627, 2002.

[11] M. Krawczyk and T. Gerkmann, "STFT phase reconstruction in voiced speech for an improved single-channel speech enhancement," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 22, no. 12, pp. 19031–1940, Dec. 2014.

[12] T. Gerkmann, M. Krawczyk, and J. Le Roux, "Phase processing for single-channel speech enhancement: History and recent advances," *IEEE Signal Process. Mag.*, vol. 32, no. 2, pp. 55–66, Mar. 2015.

[13] P. Mowlaee, R. Saeidi, and Y. Stylianou, "Advances in phase-aware signal processing in speech communication," *Speech Commun.*, vol. 81, pp. 1–29, July 2016.

[14] H. Kameoka, N. Ono, K. Kashino, and S. Sagayama, "Complex nmf: A new sparse representation for acoustic signals," in *IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, Apr. 2009, pp. 3437–3440.

[15] H. Kameoka, "Multi-resolution signal decomposition with time-domain spectrogram factorization," in *IEEE Int. Conf.*

[16] V. Emiya, R. Hamon, and C. Chaux, "Being low-rank in the time-frequency plane," in *IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, Apr. 2018, pp. 4659–4663.

[17] I. Bayram and M. E. Kamasak, "A simple prior for audio signals," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 21, no. 6, pp. 1190–1200, June 2013.

[18] P. Magron, R. Badeau, and B. David, "Phase reconstruction of spectrograms with linear unwrapping: Application to audio signal restoration," in *Eur. Signal Process.Conf. (EUSIPCO)*, Aug. 2015.

[19] K. Yatabe and Y. Oikawa, "Phase corrected total variation for audio signals," in *IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, Apr. 2018, pp. 656–660.

[20] P. Magron, R. Badeau, and B. David, "Model-based STFT phase recovery for audio source separation," *IEEE/ACM Trans. on Audio, Speech, Lang. Process.*, vol. 26, no. 6, pp. 1095–1105, June 2018.

[21] Y. Masuyama, K. Yatabe, and Y. Oikawa, "Model-based phase recovery of spectrograms via optimization on Riemannian manifolds," in *IEEE Int. Workshop Acoust. Signal Enhance. (IWAENC)*, Sept. 2018, pp. 126–130.

[22] Y. Masuyama, K. Yatabe, and Y. Oikawa, "Phase-aware harmonic/percussive source separation via convex optimization," in *IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, 2019.

[23] F. Auger and P. Flandrin, "Improving the readability of time-frequency and time-scale representations by the reassignment method," *IEEE Trans. Signal Process.*, vol. 43, no. 5, pp. 1068–1089, May 1995.

[24] N. Holighaus, Z. Průša, and P. L. Søndergaard, "Reassignment and synchrosqueezing for general timefrequency filter banks, subsampling and processing," *Signal Processing*, vol. 125, pp. 1–8, 2016.

[25] S. Fenet, R. Badeau, and G. Richard, "Reassigned timefrequency representations of discrete time signals and application to the constant-q transform," *Signal Processing*, vol. 132, pp. 170–176, Mar. 2017.

[26] C. Kehling, J. Abeßer, G. Dittmar, and G. Schuller, "Automatic tablature transcription of electric guitar recordings by estimation of score- and instrument-related parameters," in *Int. Conf. Digit. Audio. Effects*, 2014.

[27] N. Parikh and S. Boyd, "Proximal algorithms," *Found. Trends Opt.*, vol. 1, no. 3, pp. 127–239, Jan. 2014.

[28] D. Kitamura, H. Saruwatari, H. Kameoka, Y. Takahashi, K. Kondo, and S. Nakamura, "Multichannel signal separation combining directional clustering and nonnegative matrix factorization with spectrogram restoration," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 23, no. 4, pp. 654–669, Apr. 2015, Dataset: http://d-kitamura.net/en/dataset_en.htm.

[29] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Found. Trends Mach. Learn.*, vol. 3, no. 1, pp. 1–122, 2010.

[30] V. Emiya, E. Vincent, N. Harlander, and V. Hohmann, "Subjective and objective quality assessment of audio source separation," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 19, no. 7, pp. 2046–2057, Sept. 2011.