

NON-GRIFFIN-LIM TYPE SIGNAL RECOVERY FROM MAGNITUDE SPECTROGRAM

Ryusei Nakatsu

Daichi Kitahara

Akira Hirabayashi

Graduate School of Information Science and Engineering, Ritsumeikan University, Shiga, Japan

ABSTRACT

Speech and audio signal processing frequently requires to recover a time-domain signal from the magnitude of a spectrogram. Conventional methods inversely transform the magnitude spectrogram with a phase spectrogram recovered by the Griffin–Lim algorithm or its accelerated versions. The short-time Fourier transform (STFT) perfectly matches this framework, while other useful spectrogram transforms, such as the constant-Q transform (CQT), do not, because their inverses cannot be computed easily. To make the best of such useful spectrogram transforms, we propose an algorithm which recovers the time-domain signal without the inverse spectrogram transforms. We formulate the signal recovery as a nonconvex optimization problem, which is difficult to solve exactly. To approximately solve the problem, we exploit a stochastic convex optimization technique. A well-organized block selection enables us both to avoid local minimums and to achieve fast convergence. Numerical experiments show the effectiveness of the proposed method for both STFT and CQT cases.

Index Terms— Spectrogram, signal recovery, Griffin–Lim algorithm, constant-Q transform, nonconvex stochastic optimization.

1. INTRODUCTION

Spectrogram is a useful expression for analysis of speech and audio signals. Many applications, including *sound source separation* [1]–[3] and *automatic music transcription* [4], [5], are performed in the spectral-domain using spectrograms rather than in the original time-domain. A spectrogram is typically defined as a complex matrix and can be divided into *magnitude* and *phase spectrograms*. Even though both components constitute the original spectrogram, only the magnitude spectrogram is manipulated in some applications [6]–[9]. For example, to extract a specific target sound from a mixed sound, the magnitude spectrogram of the target sound is estimated from that of the mixed sound by applying a well-designed time-frequency mask. Then, a time-domain signal, which is an estimate of the target sound, is recovered through the inverse transform from the estimated magnitude spectrogram with the observed mixed phase spectrogram [10].

However, the magnitude spectrogram of the recovered signal using the mixed phase will be different from the estimated magnitude. To bring the magnitude spectrogram of the recovered signal closer to the estimated one, we have to find a more appropriate phase spectrogram. This problem is often called *phase recovery* or *phase retrieval*. In general, the size of a sound signal is relatively large, and thus the *PhaseLift* approach [11]–[13] is difficult to use since it requires huge memory resources. A well-known method for the phase recovery of the sound signal is the Griffin–Lim algorithm (GLA) [14], which performs *alternating projections*. Since its convergence is slow, accelerated versions are proposed [15], [16]. Particularly, the method in [16] uses the *alternating direction method of multipliers (ADMM)* [17].

This work was supported in part by JSPS Grants-in-Aid Grant Number for Early-Career Scientists JP19K20361 (e-mail: d-kita@fc.ritsumei.ac.jp).

In the above Griffin–Lim (GL) type approach [14]–[16], the inverse of the spectrogram transform has to be computed both for update of the phase spectrogram and for the post-processing to recover a time-domain signal from the estimated complex spectrogram. The short-time Fourier transform (STFT), the most commonly used spectrogram transform, perfectly matches this framework because its inverse is computed easily thanks to the nature of the discrete Fourier transform. However, other useful spectrogram transforms such as the constant-Q transform (CQT) [18], [19] and wavelet transforms [20], which are often used for analysis of overtone components, are not fit for this framework because their inverses cannot be computed easily.

To make the best of such useful spectrogram transforms, we propose an algorithm which can directly recover a time-domain signal from the magnitude spectrogram without computing the inverse of the spectrogram transform. To this end, we first formulate the time-domain signal recovery using the time-domain signal itself while GL type algorithms [14]–[16] used a complex spectrogram as a variable to describe the problem. The formulated problem is nonconvex and difficult to solve exactly. Instead, we approximately solve the problem by applying a stochastic convex optimization technique to avoid local minimums and achieve fast convergence. Specifically, we use the *stochastic dual coordinate ascent ADMM (SDCA-ADMM)* [21], which exploits an inexact augmented Lagrangian and thus can avoid the inverse computation for any spectrogram transform. In the proposed method, the spectrogram is divided into several blocks because of the stochastic optimization. By randomly choosing *multiple blocks* in a well-organized novel manner, the proposed method can recover the time-domain signal with the same convergence performance as in [16]. Numerical experiments show the effectiveness of the proposed non-GL type signal recovery for STFT and CQT spectrograms.

2. GRIFFIN-LIM TYPE PHASE RECOVERY

Let \mathbb{R} , \mathbb{R}_+ , and \mathbb{C} be the sets of all real, nonnegative real, and complex numbers, respectively. The imaginary unit is denoted by $i \in \mathbb{C}$. For any $c \in \mathbb{C}$, \bar{c} , $|c| \in \mathbb{R}_+$, and $\angle c \in (-\pi, \pi]$ stand for the complex conjugate, magnitude, and phase of c , respectively. We write vectors with boldface small letters and matrices with capital letters. We denote the transpose, adjoint, and inverse operators by $(\cdot)^T$, $(\cdot)^H$, and $(\cdot)^{-1}$, respectively. The composition of mappings and the Hadamard product are denoted by \circ and \odot . The floor function is denoted by $\lfloor \cdot \rfloor$.

2.1. Griffin–Lim Algorithm

Given a discrete-time signal $x[t]$, let $X := (\hat{x}_{i,j}) \in \mathbb{C}_+^{L \times J}$ be the spectrogram defined by the short-time Fourier transform (STFT), as

$$\hat{x}_{i,j} := \sum_{\tau=0}^{L-1} w[\tau] x[(j-1)\xi + \tau] e^{-i \frac{2\pi(i-1)\tau}{L}},$$

where L is the frame length, J is the number of frames, $\xi (< L)$ is the shift of a window function $w[\tau]$, and $\mathbb{C}_+^{L \times J} := \{(\hat{x}_{i,j}) \in \mathbb{C}^{L \times J} \mid \hat{x}_{1,j} \in \mathbb{R} (\forall j) \text{ and } \hat{x}_{i,j} = \hat{x}_{L-i+2,j} (\forall j \text{ and } i = 2, 3, \dots, \lfloor \frac{L+2}{2} \rfloor)\}$

stands for the *spectrogram space*. We denote STFT by a linear mapping $\mathcal{S} : \mathbb{R}^{(J-1)\xi+L} \ni \mathbf{x} \mapsto \mathcal{S}(\mathbf{x}) = X \in \mathbb{C}_{\downarrow}^{L \times J}$. Note that, in this definition, the spectrogram X includes $\hat{x}_{i,j}$ s.t. $i \in [\lfloor \frac{L+2}{2} \rfloor + 1, L]$.

Since $\xi < L$, the dimension JL of the spectrogram space $\mathbb{C}_{\downarrow}^{L \times J}$ is greater than that of the time signal space given by $(J-1)\xi + L$. Therefore the *range* of STFT $\mathcal{R}_{\mathcal{S}} := \{X \in \mathbb{C}_{\downarrow}^{L \times J} \mid \exists \mathbf{x} X = \mathcal{S}(\mathbf{x})\}$ is a low-dimensional subspace of $\mathbb{C}_{\downarrow}^{L \times J}$. By exploiting this property, GLA finds a spectrogram $X \in \mathbb{C}_{\downarrow}^{L \times J}$ which has an appropriate phase for a given magnitude spectrogram $M \in \mathbb{R}_{+\downarrow}^{L \times J} (\subset \mathbb{C}_{\downarrow}^{L \times J})$ in the intersection of $\mathcal{R}_{\mathcal{S}}$ and $\mathcal{M} := \{X \in \mathbb{C}_{\downarrow}^{L \times J} \mid |X| := (\|\hat{x}_{i,j}\|) = M\}$. More specifically, GLA solves the following feasibility problem [14]:

$$\text{find } X \in \mathcal{M} \cap \mathcal{R}_{\mathcal{S}} \quad (1)$$

by alternating projections

$$X^{(l+1)} = \mathcal{P}_{\mathcal{R}_{\mathcal{S}}}(\mathcal{P}_{\mathcal{M}}(X^{(l)})), \quad (2)$$

where $\mathcal{P}_{\mathcal{M}}$ is the projection onto \mathcal{M} given by

$$\mathcal{P}_{\mathcal{M}}(X) = M \odot \Theta_X \quad (3)$$

with the phase matrix $\Theta_X := (e^{i\angle \hat{x}_{i,j}}) \in \mathbb{C}_{\downarrow}^{L \times J}$ of X , and $\mathcal{P}_{\mathcal{R}_{\mathcal{S}}}$ is the projection onto $\mathcal{R}_{\mathcal{S}}$ given by

$$\mathcal{P}_{\mathcal{R}_{\mathcal{S}}}(X) = \mathcal{S} \circ (\mathcal{S}^H \circ \mathcal{S})^{-1} \circ \mathcal{S}^H(X). \quad (4)$$

The computation of $\mathcal{P}_{\mathcal{M}}$ in (3) is very easy since we only have to replace the magnitude of X with M . In the right-hand side of (4), $(\mathcal{S}^H \circ \mathcal{S})^{-1} \circ \mathcal{S}^H$ is equivalent to the inverse STFT (ISTFT). Note that $\mathcal{S}^H \circ \mathcal{S}$ for STFT becomes a diagonal matrix, and its inverse is computed simply by the reciprocal of each diagonal component. The corresponding matrix for CQT does not, however, become diagonal. Hence, GLA is executable for STFT, but not for CQT (see Sect. 2.3).

2.2. Griffin–Lim Like Phase Recovery via ADMM

GLA assumed that $\mathcal{M} \cap \mathcal{R}_{\mathcal{S}} \neq \emptyset$ on the problem of (1). In general, the estimated magnitude spectrogram M includes some errors, and thus $\mathcal{M} \cap \mathcal{R}_{\mathcal{S}} \neq \emptyset$ is not guaranteed. To obtain an accurate phase spectrogram even if $\mathcal{M} \cap \mathcal{R}_{\mathcal{S}} = \emptyset$, Masuyama *et al.* proposed an algorithm which finds the spectrogram closest to $\mathcal{R}_{\mathcal{S}}$ among spectrograms belonging to \mathcal{M} by the following optimization problem [16]:

$$\underset{X \in \mathcal{M}}{\text{minimize}} \quad \frac{1}{2} \|X - \mathcal{P}_{\mathcal{R}_{\mathcal{S}}}(X)\|_{\text{F}}^2, \quad (5)$$

where $\|\cdot\|_{\text{F}}$ denotes the Frobenius norm. To solve the problem of (5) faster than GLA, they adopted ADMM [17], which is briefly summarized as follows. For two proper lower-semicontinuous convex functions $f : \mathbb{R}^m \rightarrow \mathbb{R} \cup \{\infty\}$ and $g : \mathbb{R}^k \rightarrow \mathbb{R} \cup \{\infty\}$, let us

$$\underset{\mathbf{x} \in \mathbb{R}^m, \mathbf{z} \in \mathbb{R}^k}{\text{minimize}} \quad f(\mathbf{x}) + g(\mathbf{z}) \quad \text{subject to } A\mathbf{x} + B\mathbf{z} = \mathbf{c}, \quad (6)$$

where $A \in \mathbb{R}^{n \times m}$, $B \in \mathbb{R}^{n \times k}$, and $\mathbf{c} \in \mathbb{R}^n$. For the problem of (6), an augmented Lagrangian function is constructed as

$$\begin{aligned} \mathcal{L}_{\rho}(\mathbf{x}, \mathbf{z}, \mathbf{u}) := & f(\mathbf{x}) + g(\mathbf{z}) - \mathbf{u}^T (A\mathbf{x} + B\mathbf{z} - \mathbf{c}) \\ & + \frac{\rho}{2} \|A\mathbf{x} + B\mathbf{z} - \mathbf{c}\|_2^2, \end{aligned}$$

where $\mathbf{u} \in \mathbb{R}^n$ is the dual variable and $\rho > 0$. By minimizing \mathcal{L}_{ρ} in terms of \mathbf{x} and \mathbf{z} alternately, and by increasing \mathcal{L}_{ρ} in terms of \mathbf{u} , the following update formulas converging to the solution are given:

$$\begin{cases} \mathbf{x}^{(l+1)} = \underset{\mathbf{x} \in \mathbb{R}^m}{\text{argmin}} \quad f(\mathbf{x}) + \frac{\rho}{2} \|A\mathbf{x} + B\mathbf{z}^{(l)} - \mathbf{c} - \frac{1}{\rho} \mathbf{u}^{(l)}\|_2^2, \\ \mathbf{z}^{(l+1)} = \underset{\mathbf{z} \in \mathbb{R}^k}{\text{argmin}} \quad g(\mathbf{z}) + \frac{\rho}{2} \|A\mathbf{x}^{(l+1)} + B\mathbf{z} - \mathbf{c} - \frac{1}{\rho} \mathbf{u}^{(l)}\|_2^2, \\ \mathbf{u}^{(l+1)} = \mathbf{u}^{(l)} - \rho(A\mathbf{x}^{(l+1)} + B\mathbf{z}^{(l+1)} - \mathbf{c}). \end{cases} \quad (7)$$

When using ADMM, it is important to check whether or not argmin in the first and second lines of (7) can be computed easily.

Although the problem of (5) is nonconvex due to \mathcal{M} , Masuyama *et al.* applied (7) and proposed the following GL like phase recovery:

$$\begin{cases} X^{(l+1)} = \mathcal{P}_{\mathcal{M}}(Z^{(l)} + \frac{1}{\rho} U^{(l)}), \\ Z^{(l+1)} = \frac{1}{\rho+1} [\rho X^{(l+1)} - U^{(l)} + \mathcal{P}_{\mathcal{R}_{\mathcal{S}}}(X^{(l+1)} - \frac{1}{\rho} U^{(l)})], \\ U^{(l+1)} = U^{(l)} - \rho(X^{(l+1)} - Z^{(l+1)}), \end{cases} \quad (8)$$

by defining \mathbf{x} in (6) as $X \in \mathbb{C}_{\downarrow}^{L \times J}$, \mathbf{z} as $Z \in \mathbb{C}_{\downarrow}^{L \times J}$, f as the indicator function of \mathcal{M} , g as the cost function in (5), A as the identity matrix, B as the negative identity matrix, and \mathbf{c} as the zero matrix. This algorithm achieved better convergence performance than GLA, but $\mathcal{P}_{\mathcal{R}_{\mathcal{S}}}$ appears again in the second line. Thus this is not executable for CQT.

2.3. Constant-Q Transform

The constant-Q transform (CQT) [18], [19] is useful for analysis of overtone components because it analyzes frequency components in the following geometric progression with the common ratio $2^{1/r}$:

$$f_i := 2^{\frac{i-1}{r}} f_{\min} > 0 \quad (i = 1, 2, \dots, I_{\max}),$$

where $f_{\max} := f_{I_{\max}} < f_s/2$, f_s is the sampling frequency, and r is often set to 12 or 24. The name ‘CQT’ is from the fact a quality factor

$$Q_i := \left\lfloor q \frac{f_i}{\delta f_i} \right\rfloor := \left\lfloor q \frac{f_i}{f_{i+1} - f_i} \right\rfloor = \left\lfloor q \frac{1}{2^{1/r} - 1} \right\rfloor = Q$$

is *constant* independently of i , where $q \in (0, 1]$ and $\lfloor \cdot \rfloor$ is the round-down function. For example, if $r = 24$ and $q = 0.82$, then $Q = 28$.

Components of a CQT spectrogram $X := (\hat{x}_{i,j})$ are defined by

$$\hat{x}_{i,j} := \frac{1}{\sqrt{L_i}} \sum_{\tau=0}^{L_i-1} w_i[\tau] x[(j-1)\xi_i + \tau] e^{-i \frac{2\pi Q \tau}{L_i}},$$

for $i = 1, 2, \dots, I_{\max}$, where $L_i := \lfloor Q f_s / f_i \rfloor$ is the frame length and $\xi_i (< L_i)$ is the shift of a window function $w_i[\tau]$ for the i th frequency. We also define $\hat{x}_{2I_{\max}-i+1,j} := \hat{x}_{i,j}$ for $i = 1, 2, \dots, I_{\max}$ as components corresponding to $-f_i$. Note that a CQT spectrogram X is a matrix only if $\forall i \xi_i = \xi < L_{I_{\max}}$ holds, and in this case the spectrogram space is defined as $\mathbb{C}_{\downarrow}^{2I_{\max} \times J} := \{(\hat{x}_{i,j}) \in \mathbb{C}^{2I_{\max} \times J} \mid \hat{x}_{i,j} = \hat{x}_{2I_{\max}-i+1,j} (\forall j \text{ and } i = 1, 2, \dots, I_{\max})\}$, where J is the number of frames. If we denote CQT by a linear mapping \mathcal{S} , $\mathcal{S}^H \circ \mathcal{S}$ is not a diagonal matrix and ICQT $(\mathcal{S}^H \circ \mathcal{S})^{-1} \circ \mathcal{S}^H$ is not available.

3. NON-GRIFFIN-LIM TYPE SIGNAL RECOVERY WITHOUT THE INVERSE SPECTROGRAM TRANSFORM

To sum up the discussion in Sect. 2, the GL type *phase recovery* algorithms [14]–[16] require the computation of $(\mathcal{S}^H \circ \mathcal{S})^{-1}$, which is easy for STFT because $\mathcal{S}^H \circ \mathcal{S}$ is a simple diagonal matrix. However, when \mathcal{S} stands for CQT, it is difficult to compute $(\mathcal{S}^H \circ \mathcal{S})^{-1}$, and to apply the update formulas of (2) and (8). In addition, even if the solution \hat{X} to the optimization problem such as (1) or (5) can be obtained, the inverse transform $(\mathcal{S}^H \circ \mathcal{S})^{-1} \circ \mathcal{S}^H$ has to be computed as the post-processing to recover a time-domain signal \mathbf{x} from \hat{X} .

To eliminate the inverse transform $(\mathcal{S}^H \circ \mathcal{S})^{-1} \circ \mathcal{S}^H$ as the post-processing, we reformulate the optimization problem of (5) as an essentially equivalent problem using the time-domain signal itself:

$$\underset{\mathbf{x} \in \mathbb{R}^{(J-1)\xi+L}}{\text{minimize}} \quad \frac{1}{2} \|M - |\mathcal{S}(\mathbf{x})|\|_{\mathbb{F}}^2. \quad (9)$$

We propose a non-GL type *signal recovery* algorithm which reconstructs \mathbf{x} by solving (9) without computing $(\mathcal{S}^H \circ \mathcal{S})^{-1}$. To this end, we use a stochastic convex optimization technique SDCA-ADMM.

3.1. Stochastic Dual Coordinate Ascent ADMM

SDCA-ADMM exploits the characteristic that the dual variable \mathbf{u} in (7) converges to the optimal solution of the *dual problem* [21]. The convex optimization problem to be considered is to

$$\underset{\mathbf{x} \in \mathbb{R}^m}{\text{minimize}} \quad \sum_{i=1}^n f_i(\mathbf{a}_i^T \mathbf{x}) + g(B^T \mathbf{x}), \quad (10)$$

where $A := (\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n) \in \mathbb{R}^{m \times n}$, $B \in \mathbb{R}^{m \times h}$, and $f: \mathbb{R}^n \ni \mathbf{u} \mapsto \sum_{i=1}^n f_i(u_i) \in \mathbb{R} \cup \{\infty\}$ and $g: \mathbb{R}^h \rightarrow \mathbb{R} \cup \{\infty\}$ are supposed to be a data fidelity and a regularization term, respectively. The dual problem for the primal problem of (10) is to

$$\underset{\mathbf{u} \in \mathbb{R}^n, \mathbf{v} \in \mathbb{R}^h}{\text{minimize}} \quad \sum_{i=1}^n f_i^*(u_i) + g^*(\mathbf{v}) \quad \text{subject to } A\mathbf{u} + B\mathbf{v} = \mathbf{0}, \quad (11)$$

where f_i^* and g^* are the *conjugate function* of f_i and g , respectively. By solving (11) with the ADMM update formulas as in (7), the *dual* variable \mathbf{u} converges to the solution to the problem of (10), because the relationship between \mathbf{x} and \mathbf{u} in (7) is reversed.

In many cases, however, it is difficult to directly adopt the above idea because argmin in the first and second lines of (7) cannot be computed easily. To overcome the difficulty, SDCA-ADMM uses an *inexact* augmented Lagrangian function, which is defined by adding two proximal terms as

$$\begin{aligned} \widehat{\mathcal{L}}_\rho^{(l)}(\mathbf{u}, \mathbf{v}, \mathbf{x}) &:= f^*(\mathbf{u}) + g^*(\mathbf{v}) - \mathbf{x}^T (A\mathbf{u} + B\mathbf{v}) \\ &+ \frac{\rho}{2} \|A\mathbf{u} + B\mathbf{v}\|_2^2 + \frac{1}{2} \|\mathbf{u} - \mathbf{u}^{(l)}\|_{G_A}^2 + \frac{1}{2} \|\mathbf{v} - \mathbf{v}^{(l)}\|_{G_B}^2, \end{aligned}$$

where $G_A := \rho(\eta_A I_n - A^T A)$ and $G_B := \rho(\eta_B I_h - B^T B)$ are positive semidefinite matrices to avoid the computation of $(A^T A)^{-1}$ and $(B^T B)^{-1}$. SDCA-ADMM minimizes $\widehat{\mathcal{L}}_\rho^{(l)}$ in terms of \mathbf{u} and \mathbf{v} alternately and increases it in terms of \mathbf{x} . Note that the first term in (11) is separated in terms of u_i so that the stochastic optimization can be easily introduced as follows. SDCA-ADMM splits the index set $\mathcal{I} := \{1, 2, \dots, n\}$ for f_i^* into b subsets $\{\mathcal{I}_d\}_{d=1}^b$ such that $\mathcal{I}_d \neq \emptyset$, $\bigcup_{d=1}^b \mathcal{I}_d = \mathcal{I}$ and $\mathcal{I}_d \cap \mathcal{I}_{d'} = \emptyset$ ($d \neq d'$). Then, SDCA-ADMM updates the variables by

$$\left[\begin{array}{l} \text{Choose one index } d \in \{1, 2, \dots, b\} \text{ uniformly at random,} \\ \mathbf{q}^{(l+1)} = \mathbf{v}^{(l)} + \frac{1}{\rho\eta_B} B^T [\mathbf{x}^{(l)} - \rho(A\mathbf{u}^{(l)} + B\mathbf{v}^{(l)})], \\ \mathbf{v}^{(l+1)} = \mathbf{q}^{(l+1)} - \frac{1}{\rho\eta_B} \text{prox}_{\rho\eta_B g}(\rho\eta_B \mathbf{q}^{(l+1)}), \\ \mathbf{p}_{\mathcal{I}_d}^{(l+1)} = \mathbf{u}_{\mathcal{I}_d}^{(l)} + \frac{1}{\rho\eta_{A_{\mathcal{I}_d}}} A_{\mathcal{I}_d}^T [\mathbf{x}^{(l)} - \rho(A\mathbf{u}^{(l)} + B\mathbf{v}^{(l+1)})], \\ \mathbf{u}_{\mathcal{I}_d}^{(l+1)} = \mathbf{p}_{\mathcal{I}_d}^{(l+1)} - \frac{1}{\rho\eta_{A_{\mathcal{I}_d}}} \text{prox}_{\rho\eta_{A_{\mathcal{I}_d}} f_{\mathcal{I}_d}}(\rho\eta_{A_{\mathcal{I}_d}} \mathbf{p}_{\mathcal{I}_d}^{(l+1)}), \\ \mathbf{x}^{(l+1)} = \mathbf{x}^{(l)} - \gamma\rho [A\mathbf{u}^{(l+1)} + B\mathbf{v}^{(l+1)} \\ \quad - (1 - \phi)(A\mathbf{u}^{(l)} + B\mathbf{v}^{(l)})], \end{array} \right. \quad (12)$$

where $\gamma > 0$, $\eta_{A_{\mathcal{I}_d}} \geq \|A_{\mathcal{I}_d}\|_{\text{op}}^2$, $\eta_B \geq \|B\|_{\text{op}}^2$, $\|\cdot\|_{\text{op}}$ is the operator

norm, $\phi = 1/b$ is the probability that each u_i is chosen, and prox denotes the *proximity operator*. When the proximity operators of f_i and g are computed easily, the update formulas of (12) can be executed.

3.2. Signal Recovery via SDCA-ADMM

We solve the nonconvex optimization problem of (9) via the SDCA-ADMM update formulas. Define $f(\mathbf{u}) = \sum_i f_i(u_i)$ in (10) as

$$\sum_{i,j} f_{i,j}(\hat{u}_{i,j}) := \sum_{i,j} \frac{1}{2} (m_{i,j} - |\hat{u}_{i,j}|)^2 = \frac{1}{2} \|M - |U|\|_{\mathbb{F}}^2,$$

g as the zero function, A as \mathcal{S}^H , and B as the zero matrix. Then, the dual problem of (9) is given by

$$\underset{U \in \mathbb{C}^{L \times J}}{\text{minimize}} \quad \sum_{i,j} f_{i,j}^*(\hat{u}_{i,j}) \quad \text{subject to } \mathcal{S}^H(U) = \mathbf{0}. \quad (13)$$

To solve the problem of (13) using SDCA-ADMM, we have to carefully choose the index set $\mathcal{I}_d \subset \mathcal{I} := \{(i, j)\}_{j=1,2,\dots,L}^{i=1,2,\dots,L}$. In the conventional method, only one index set out of the b index sets was randomly chosen as in (12). This strategy not only prevents acceleration but also may get stuck in a local minimum due to the nonconvexity.

To overcome this difficulty, it should be noted that the convergence of (12) is guaranteed if the parameter ϕ agrees with the probability that each component u_i is chosen from \mathbf{u} . Therefore, instead of updating one of the b non-overlapping index sets, we propose to update *multiple index sets* with a probability $\phi \in (0, 1]$, where $\phi = 1$ means the non-stochastic optimization. Based on these discussions, we propose the following update formulas for solving (13):

$$\left[\begin{array}{l} \text{Choose multiple indices } \mathcal{D} \in 2^{\{1,2,\dots,b\}} \text{ with a probability } \phi, \\ P_{\mathcal{I}_D}^{(l+1)} = U_{\mathcal{I}_D}^{(l)} + \frac{1}{\rho\eta} \mathcal{S}_{\mathcal{I}_D}(\mathbf{x}^{(l)} - \rho\mathcal{S}^H(U^{(l)})), \\ U_{\mathcal{I}_D}^{(l+1)} = \frac{1}{\rho\eta+1} (\rho\eta P_{\mathcal{I}_D}^{(l+1)} - M_{\mathcal{I}_D} \odot \Theta_{P_{\mathcal{I}_D}^{(l+1)}}), \\ \mathbf{x}^{(l+1)} = \mathbf{x}^{(l)} - \gamma\rho\mathcal{S}^H(U^{(l+1)} - (1 - \phi)U^{(l)}), \end{array} \right. \quad (14)$$

where $\rho > 0$, $\gamma > 0$, and $\eta \geq \|\mathcal{S}\|_{\text{op}}^2$. Note that there is no computation of the inverse spectrogram transform $(\mathcal{S}^H \circ \mathcal{S})^{-1} \circ \mathcal{S}^H$ in the formulas of (14). Hence, the proposed algorithm is *non-GL type* and can be applied not only to STFT but also to CQT. As shown in Sect. 4.2, we propose to divide the index set along the frequency axis, as $\{(1, j)\}_{j=1}^J$ and $\{(i, j), (L-i+2, j)\}_{j=1}^J$ ($i = 2, 3, \dots, \lfloor \frac{L+2}{2} \rfloor$).

4. NUMERICAL EXPERIMENTS

To confirm the convergence performance of the proposed algorithm and the validity of the stochastic optimization, we conducted numerical experiments of time-domain signal recovery from a STFT magnitude spectrogram in Sects. 4.1 & 4.2 and from a CQT magnitude spectrogram in Sect. 4.3. Each recovered time-domain signal \mathbf{x} was evaluated by the normalized error $100 \frac{\|M - |\mathcal{S}(\mathbf{x})|\|_{\mathbb{F}}}{\|M\|_{\mathbb{F}}} [\%]$ of the magnitude spectrogram because it corresponds to the cost function in (9) and phase errors have less effect on human ears. An open source piano sound of 7,000 ms and the Hanning window were used with sampling frequency 16 kHz, frame length 64 ms, and frame shift 32 ms. In all methods (2), (8) and (14), the number of iterations was 10,000.

4.1. Performance Evaluation of the Proposed Method in STFT

We compared the performance of the proposed method and the existing methods [14], [16] under the ideal situation of $\mathcal{M} \cap \mathcal{R}_S \neq \emptyset$.

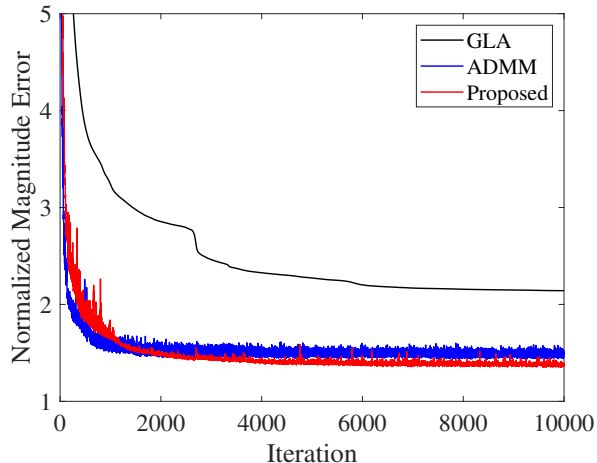


Fig. 1. Normalized error curves for 3 methods (STFT case).

The piano sound was recovered by each method and the normalized error was computed. The parameter in the ADMM method [16] was $\rho = 0.08$ and those in the proposed method were $\gamma = 1$, $\rho = 0.01$, $\eta = 1.001L$, and $\phi = 0.75$. The index set $\mathcal{I} = \{(i, j)\}_{j=1,2,\dots,L}^{i=1,2,\dots,L}$ was divided along the frequency axis as mentioned in Sect. 3.2. The parameter values were empirically adjusted for the best results. Figure 1 shows the normalized error curve in terms of the iteration for each method. We can see that the ADMM method and the proposed method indicate faster convergence speed than GLA. We emphasize that the proposed method achieves mostly the same convergence performance as the ADMM method *without ISTFT* $(\mathcal{S}^H \circ \mathcal{S})^{-1} \circ \mathcal{S}^H$.

4.2. Validity Evaluation of the Stochastic Optimization in STFT

We prepared three types of divisions of the index set \mathcal{I} : (a) *primal division* $\{(1, j)\}_{j=1,2,\dots,J}$ and $\{(i, j), (L-i+2, j)\}_{i=2,3,\dots,\lfloor \frac{L+2}{2} \rfloor; j=1,2,\dots,J}$ ($b = J \lfloor \frac{L+2}{2} \rfloor$), (b) *frame division* $\{(i, j)\}_{i=1}^L$ ($j = 1, 2, \dots, J$) ($b = J$), and (c) *frequency division* $\{(1, j)\}_{j=1}^J$ and $\{(i, j), (L-i+2, j)\}_{j=1}^J$ ($i = 2, 3, \dots, \lfloor \frac{L+2}{2} \rfloor$) ($b = \lfloor \frac{L+2}{2} \rfloor$). We also prepared the selection probability ϕ for each division as $\phi = 0.5$ and $\phi = 0.75$. In addition, $\phi = 1$ was also used for comparison to the non-stochastic optimization scenario. Figure 2 shows the normalized error curves in terms of the iteration.¹ All the curves except for that of the frame division with $\phi = 0.75$ are lower than that for $\phi = 1$. Therefore, we can confirm the effectiveness of the stochastic optimization. The best performance was given by the frequency division with $\phi = 0.75$.

4.3. Validity Evaluation of the Stochastic Optimization in CQT

We also conducted numerical experiments of time-domain signal recovery from a CQT magnitude spectrogram by the proposed method. We set the parameters of CQT to $r = 24$, $q = 0.82$, $f_{\min} = 60$ Hz, $f_{\max} = 5,922$ Hz (i.e., $I_{\max} = 160$), and $\xi_i = \xi = L_{I_{\max}}/2$. Since we cannot use the frame division if $\xi_i \neq \xi_{i'}$ ($i \neq i'$), we used only the primal division and the frequency division. Figure 3 shows the normalized magnitude error curves in terms of the iteration for the proposed method. The lowest normalized error converged to 0.63 %, which means that the proposed method enables time-domain signal

¹Strictly speaking, the ten point moving average filter is used to make the curves more visible.

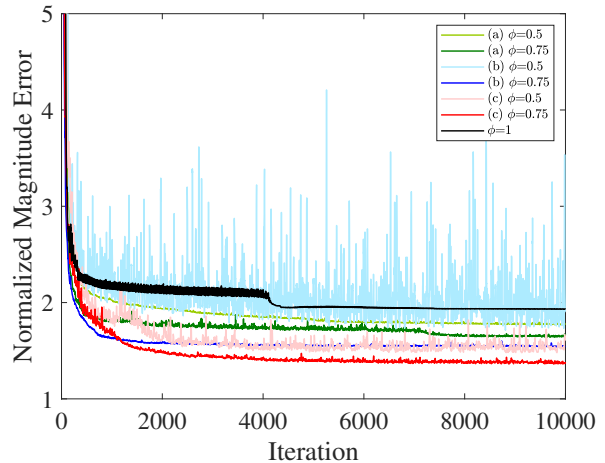


Fig. 2. Normalized error curves for 3 index set divisions (STFT case).

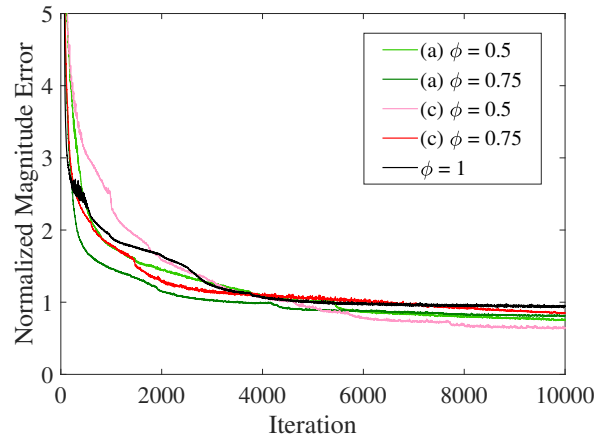


Fig. 3. Normalized error curves for 2 index set divisions (CQT case).

recovery from the CQT magnitude spectrogram. In each type of division of the index set \mathcal{I} , there was a trade-off between the convergence speed and the normalized error according to the probability ϕ . When $\phi = 0.5$, the normalized error was lower than when $\phi = 0.75$. On the other hand, when $\phi = 0.75$, the convergence was faster than when $\phi = 0.5$. Since the proposed method makes the stochastic optimization more flexible, it would be best to use relatively large ϕ in the early part and then to use relatively small ϕ in the latter part.

5. CONCLUSION

In this paper, we proposed a time-domain signal recovery algorithm which does not require the computation of the inverse of the spectrogram transform. We first reformulated the signal recovery problem as a nonconvex optimization problem using the time-domain signal itself. To solve this problem, we used a stochastic optimization technique SDCA-ADMM for fast convergence without the inverse transform. Numerical experiments showed that the proposed method with a well-organized division of a cost function along the frequency axis achieves the same excellent convergence performance as the GL type algorithms for a STFT magnitude spectrogram. We also showed that the proposed method is effective for a CQT magnitude spectrogram.

6. REFERENCES

- [1] S. Makino, S. Araki, R. Mukai, and H. Sawada, "Audio source separation based on independent component analysis," in *Proceedings of IEEE International Symposium on Circuits and Systems (ISCAS)*, Vancouver, Canada, 2004, vol. 5, pp. 668–671.
- [2] T. Kim, H. T. Attias, S. Y. Lee, and T. W. Lee, "Blind source separation exploiting higher-order frequency dependencies," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 15, no. 1, pp. 70–79, 2007.
- [3] 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 Transactions on Audio, Speech, and Language Processing*, vol. 24, no. 9, pp. 1626–1641, 2016.
- [4] E. Benetos, S. Dixon, D. Giannoulis, H. Kirchhoff, and A. Klappuri, "Automatic music transcription: challenges and future directions," *Journal of Intelligent Information Systems*, vol. 41, no. 3, pp. 407–434, 2013.
- [5] M. Mauch, C. Cannam, R. Bittner, G. Fazekas, J. Salamon, J. Dai, J. Bello, and S. Dixon, "Computer-aided melody note transcription using the Tony software: accuracy and efficiency," in *Proceedings of International Conference on Technologies for Music Notation and Representation (TENOR)*, Paris, France, 2015, pp. 28–30.
- [6] G. Yu, S. Mallat, and E. Bacry, "Audio denoising by time-frequency block thresholding," *IEEE Transactions on Signal Processing*, vol. 56, no. 5, pp. 1830–1839, 2008.
- [7] J. L. Flanagan and R. M. Golden, "Phase vocoder," *Bell System Technical Journal*, vol. 45, no. 9, pp. 1493–1509, 1966.
- [8] M. Portnoff, "Implementation of the digital phase vocoder using the fast Fourier transform," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 24, no. 3, pp. 243–248, 1976.
- [9] V. Verfaillie, U. Zolzer, and D. Arfib, "Adaptive digital audio effects (A-DAFx): A new class of sound transformations," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 14, no. 5, pp. 1817–1831, 2006.
- [10] A. Jansson, E. Humphrey, N. Montecchio, R. Bittner, A. Kumar, and T. Weyde, "Singing voice separation with deep U-Net convolutional networks," in *Proceedings of International Society for Music Information Retrieval Conference (ISMIR)*, Suzhou, China, 2017, pp. 745–751.
- [11] E. J. Candès, T. Strohmer, and V. Voroninski, "PhaseLift: Exact and stable signal recovery from magnitude measurements via convex programming," *Communications on Pure and Applied Mathematics*, vol. 66, no. 8, pp. 1241–1274, 2012.
- [12] K. Jaganathan, Y. C. Eldar, and B. Hassibi, "STFT phase retrieval: Uniqueness guarantees and recovery algorithms," *IEEE Journal of Selected Topics in Signal Processing*, vol. 10, no. 4, pp. 770–781, 2016.
- [13] E. J. Candès, X. Li, and M. Soltanolkotabi, "Phase retrieval from coded diffraction patterns," *Applied and Computational Harmonic Analysis*, vol. 39, no. 2, pp. 277–299, 2015.
- [14] D. W. Griffin and J. S. Lim, "Signal estimation from modified short-time Fourier transform," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 32, no. 2, pp. 236–243, 1984.
- [15] N. Perraudin, P. Balazs, and P. L. Søndergaard, "A fast Griffin–Lim algorithm," in *Proceedings of IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, New Paltz, NY, USA, 2013, 4 pages.
- [16] Y. Masuyama, K. Yatabe, and Y. Oikawa, "Griffin–Lim like phase recovery via alternating direction method of multipliers," *IEEE Signal Processing Letters*, vol. 26, no. 1, pp. 184–188, 2019.
- [17] D. Gabay and B. Mercier, "A dual algorithm for the solution of nonlinear variational problems via finite element approximation," *Computers & Mathematics with Applications*, vol. 2, no. 1, pp. 17–40, 1976.
- [18] J. C. Brown, "Calculation of a constant-Q spectral transform," *The Journal of the Acoustical Society of America*, vol. 89, no. 1, pp. 425–434, 1991.
- [19] A. Nagathil and R. Martin, "Evaluation of spectral transforms for music signal analysis," in *Proceedings of IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, New Paltz, NY, USA, 2013, 4 pages.
- [20] I. Daubechies, *Ten Lectures on Wavelets*, ser. CBMS-NSF Regional Conference Series in Applied Mathematics. Philadelphia, PA, USA: SIAM, 1992, vol. 61.
- [21] T. Suzuki, "Stochastic dual coordinate ascent with alternating direction multiplier method," in *Proceedings of International Conference on Machine Learning (ICML)*, Beijing, China, 2014, 9 pages.