

Towards High Quality Real-Time Signal Reconstruction from STFT Magnitude

Zdeněk Průša

Abstract—An efficient algorithm for real-time signal reconstruction from the magnitude of its short-time Fourier transform (STFT) is presented. The algorithm combines strengths of two previously published algorithms and the reconstructed signals exhibit excellent perceptual quality. We present an extensive comparison with the state-of-the-art algorithms showing that the proposed method outperforms others by far in settings capable of producing high quality signals.

Index Terms—Time-frequency, short-time Fourier transform, phase reconstruction, real-time, STFT, spectrogram

I. INTRODUCTION

In time-frequency signal processing, it is a common practice to work only with the magnitude of the STFT of a signal. However, as soon as reconstruction is desired, considering phase information is crucial. When the magnitude is modified, it is often sufficient to reuse the unmodified phase to recover the signal [1], however some spectrogram modifications might invalidate the phase and the reconstruction procedure may therefore lead to undesired artifacts [2]. In some applications, the original phase is not available at all [3]. STFT phase retrieval algorithms alleviate these problems by allowing complete disposal of the existing phase and constructing a new valid phase from scratch taking the modified magnitude.

Unfortunately, to date, available STFT phase retrieval algorithms cannot always be expected to fulfill all the requirements at the same time. For example, some algorithms require the knowledge of the entire signal and they typically need a large number of costly iterations to produce a good result [4], [5], [6], [7]. This fact disqualifies them from being used in any real-time or interactive applications. Algorithms which are capable of processing signals in real-time i.e. in the frame-by-frame manner with bounded delay [8], [9], [10], [11], tend to produce noticeable artefacts such as “phasiness” [2], metallic ringing, echo etc. for specific classes of audio signals.

In this work, we propose a real-time phase reconstruction algorithm which outperforms the state-of-the-art algorithms by a large margin in the senses of both the objective error measure introduced later and the perceived quality of the reconstruction. We compare our method with the following algorithms, which can be considered to be the state-of-the-art: The Real-Time Iterative Spectrogram Inversion (RTISI) [12] later improved by including look-ahead frames [13], [8]

(RTISI-LA) and further modified slightly in the line of work of Gmann and Spiertz [14], [15], [16] (GSRTISI-LA). From the point of view of this letter, the crucial property of GSRTISI-LA is that it allows defining an initial estimate of the phase of the latest look-ahead frame.

In our previous work [17] we have proposed a non-iterative algorithm termed Real-Time Phase Gradient Heap Integration (RTPGHI). It is based on the phase-magnitude relationship which allows estimating the phase increments between neighboring coefficients solely from the magnitude. The algorithm requires one look-ahead frame (zero look-ahead frames version is also available) and, as it turns out, it is also a suitable candidate for providing the initial phase guess for GSRTISI-LA.

In this letter, we combine the strengths of both RTPGHI and GSRTISI-LA to perform high quality signal reconstruction from the spectrogram. To that end, in addition to the RTPGHI initialization, we further generalize GSRTISI-LA such that the analysis window, window overlap, number of frequency bins and the number of look-ahead frames can be chosen freely and independent of each other (to the extent specified next).

Because we aim for a high reconstruction quality, in our comparisons we do not include algorithms presented in [11], [9], [10]. Although they are much faster than the iterative algorithms, they simply produce results of significantly lower quality.

In the spirit of reproducible research, the implementation of the algorithms, audio examples as well as the scripts reproducing the experiments are available at <http://lftat.github.io/notes/048>. The code depends on our Matlab/GNU Octave [18] packages LTFAT [19], [20] (version 2.1.3 or above) and PHASERET (version 0.2.0 or above). Both the toolboxes are open-source and they can be obtained from <http://lftat.github.io> and <http://lftat.github.io/phaseret>, respectively.

The letter is organized as follows. In Section II we recall essential formulas for computing STFT analysis and synthesis, Section III contains description of the proposed algorithm and, finally, in Section IV we compare the proposed method with the state-of-the-art algorithms.

II. STFT AND ITS INVERSE

The discrete STFT of an input signal $f \in \ell^2(\mathbb{Z})$ using analysis window $g \in \ell^2(\mathbb{Z})$ is defined as

$$c_n(m) = (\mathcal{V}_g f)_n(m) = \sum_{l \in \mathbb{Z}} f(l + na) \overline{g(l)} e^{-i2\pi ml/M}, \quad (1)$$

where overline denotes complex conjugation, M is a finite number of frequency channels indexed as $m = 0, \dots, M-1$, $n \in \mathbb{Z}$ is a time-frame index and the parameter a is the time

Z. Průša* is with the Acoustics Research Institute, Austrian Academy of Sciences, Wohllebengasse 12–14, 1040 Vienna, Austria, email: zdenek.prusa@oeaw.ac.at (corresponding address).

Manuscript received April 19, 2005; revised August 26, 2015.

This work was supported by the Austrian Science Fund (FWF) START-project FLAME (“Frames and Linear Operators for Acoustical Modeling and Parameter Estimation”; Y 551-N13).

step (window shift) in samples. The window g will further be considered to be real, whole-point symmetric and finitely supported such that the index set of the summation can be reduced to

$$\mathcal{I} = \{-\lfloor \text{len}(g)/2 \rfloor, \dots, \lceil \text{len}(g)/2 \rceil - 1\}, \quad (2)$$

where $\text{len}(g)$ is the length of the window support and the center sample of the window is at index $l = 0$. The synthesis window \tilde{g} can be obtained as

$$\tilde{g} = \frac{1}{M} \frac{g}{\sum_{n \in \mathbb{Z}} g^2(\cdot - na)}, \quad (3)$$

if the following conditions are met: support of the window g is less or equal to the number of frequency channels i.e. $\text{len}(g) \leq M$ and there is some window overlap i.e. $\text{len}(g) > a$. Under these assumptions, the sum in the denominator is nonzero and a -periodic, \tilde{g} and g have equal time support and the following holds

$$\sum_{n \in \mathbb{Z}} (g\tilde{g})(\cdot - na) \equiv 1/M. \quad (4)$$

Please refer to [21], [22], [23], for example, for a thorough mathematical treatment of the invertibility of discrete STFT (Discrete Gabor transform) in the context of frame theory. In the terms of Gabor frame theory the window computed using (3) is the *canonical dual window* and the inequality $gl \leq M$ is referred to as the *painless* condition [24].

Having the synthesis window \tilde{g} , the individual time-frames f_n can be recovered using

$$f_n(l) = \begin{cases} \tilde{g}(l) \sum_{m=0}^{M-1} c_n(m) e^{i2\pi ml/M} & \text{for all } l \in \mathcal{I}, \\ 0 & \text{otherwise.} \end{cases} \quad (5)$$

A part of the signal recovered from frames up to time-frame N is given by

$$\tilde{f}_N = \sum_{n=-\infty}^N f_n(\cdot - na) \quad (6)$$

(cf. overlap-add procedure) and, clearly, the original signal can be recovered using $N = \infty$.

III. ALGORITHMS

In the following, we will denote the magnitude of coefficients of n -th time-frame as $s_n = |c_n|$. The goal of real-time phase reconstruction algorithms is to estimate coefficients c_n . We will denote the estimated coefficients as \tilde{c}_n and the coefficients can be plugged into (5) thus recovering the time frame.

A. Overview of RTPGHI

The RTPGHI algorithm [17] is a real-time capable version of the PGHI algorithm [25]. Both algorithms are based on the relationship of the gradients of the phase and the logarithm of the magnitude of STFT and employ an adaptive integration scheme to recover the phase. Best performance is achieved by using the Gaussian window, but other windows can be used as well. The RTPGHI algorithm comes in two versions requiring one (RTPGHI(1)) or zero (RTPGHI(0)) look-ahead frames respectively. For details please see the above mentioned references.

Algorithm 1: RTPGHI(1) + GSRTISI-LA($N_{\text{LA}} - 1$), n -th time frame

Input: Number of look-ahead frames N_{LA} , number of iterations I , magnitude of STFT coefficients $s_n, \dots, s_{n+N_{\text{LA}}}$

Output: Time frame f_n .

- 1 Compute $f_{n+N_{\text{LA}}-1}$ using (5) and coefficients $\tilde{c}_{n+N_{\text{LA}}-1}$ estimated using the RTPGHI algorithm (requires $s_{n+N_{\text{LA}}}$)
- 2 **for** $i = 1, 2, \dots, I$ **do**
- 3 **for** $p = N_{\text{LA}} - 1, \dots, 0$ **do**
- 4 Compute $\tilde{f}_{n+N_{\text{LA}}-1}$ using (6)
- 5 $t \leftarrow (\mathcal{V}_{g_p} \tilde{f}_{n+N_{\text{LA}}-1})_{n+p}$
- 6 $c_{n+p} \leftarrow s_{n+p} t / |t|$
- 7 Compute f_{n+p} using (5)
- 8 **end**
- 9 **end**

B. Generalized GSRTISI-LA with RTPGHI Initialization

The original version of GSRTISI-LA is restricted to a single type of window to be used for both the analysis and the synthesis (i.e. $\tilde{g} = g$), and, moreover, the authors seem to have only used 75% window overlap and 3 look-ahead frames. In this section, we present a variant of GSRTISI-LA which admits any analysis window and allows free choice of the window overlap length and the number of frequency channels (up to conditions presented in Sec. II). As already mentioned, our variant employs RTPGHI to estimate the initial phase.

Assuming RTPGHI(1) is used for initialization, the algorithm processes one time-frame at a time, taking into account N_{LA} future frames. The $N_{\text{LA}} - 1$ look-ahead frames are used for the basic version of the GSRTISI-LA algorithm, one additional look-ahead frame is required by RTPGHI. In addition to windows g and \tilde{g} , the algorithm requires N_{LA} additional analysis windows $g_0, \dots, g_{N_{\text{LA}}-1}$ which are obtained as

$$g_k = M \frac{g}{g_{\text{sum}}(\cdot + ka)}, \text{ where } g_{\text{sum}} = \sum_{n=-\infty}^{N_{\text{LA}}-1} (g\tilde{g})(\cdot - na). \quad (7)$$

The proposed algorithm RTPGHI(1) + GSRTISI-LA($N_{\text{LA}} - 1$) is formally written in Alg. 1. An extension to the version RTPGHI(0) + GSRTISI-LA(N_{LA}) is straightforward. Even though we operate with infinite sum limit, in practice, due to the finite support of the windows, it is sufficient to work with time frames in the range $n - N_{\text{LB}}, \dots, n + N_{\text{LA}}$, where $N_{\text{LB}} = \lceil \text{len}(g)/a \rceil - 1$ is the number of “look-back” frames.

Please note that other types of phase initialization are possible. The authors of the original version of GSRTISI-LA proposed to do simple *phase unwrapping* [16]. Another option is to employ different algorithms such as [9] in place of RTPGHI. However, in our experience neither of these two approaches brings considerable improvements over not doing the initialization at all. Often, it is even harmful to the overall performance.

C. Real-time Deadline, Delay and Computational Complexity

The worst case execution time for a single output frame must typically be less than a/f_s (time frame shift divided by the sampling rate) seconds to meet the real-time deadline restriction. This fact limits the number of iterations I which can be done, but the exact number is entirely dependent on the computing power of the device. Since the number of look-ahead frames can be varied, we will further use the number of *per-frame* iterations to be able to directly compare different settings.

The typical delay of the real-time STFT analysis-synthesis scheme is equal to the length of the window. Each look-ahead frame of the phase reconstruction algorithm increases the delay by the length of the window shift a , therefore the overall input-output delay is $(\text{len}(g) + aN_{\text{LA}})/f_s$ seconds.

IV. EXPERIMENTS

In the experiments we used the following objective error measure, previously referred to as *spectral convergence* [26]

$$\mathcal{C} = \sqrt{\frac{\sum_{n=0}^{N-1} \sum_{m=0}^{M-1} \left(s_n(m) - \left| (\mathcal{V}_g \tilde{f})_n(m) \right| \right)^2}{\sum_{n=0}^{N-1} \sum_{m=0}^{M-1} s_n(m)^2}}, \quad (8)$$

where \tilde{f} is the reconstructed signal and N denotes the total number of time frames of the finite signal. The transform \mathcal{V}_g uses identical values for the parameters (g , a and M) as the one used to obtain s . Values in decibels are obtained by $20 \log_{10} \mathcal{C}$. In the experiments, we used the SQAM database [27] which consists of 70 recordings sampled at 44.1 kHz. Only the first 10 seconds from the first channel of each sound sample was used in the evaluation.

In our experience, a substantial window overlap is necessary in order to produce results of high perceptual quality. Therefore, in our tests, we use 87.5% window overlap, which results from using time step size $a = 256$ in conjunction with the fixed number of frequency channels $M = 2048$, and the Gaussian window as the analysis window truncated at 1% of its height such that $\text{len}(g) = 2048$. Using even higher window overlap further improves the results. Please note that whenever we refer to the average error in dB we mean $20 \log_{10} \frac{1}{70} \sum_{k=1}^{70} \mathcal{C}_k$, where \mathcal{C}_k is the error of the k -th sound excerpt obtained from (8). Averaging errors that have been already converted to dB (which is occasionally done in other contributions) produces even better (lower) errors for all the algorithms.

In the real-time setting, there is room only for a limited number of iterations, but since the exact number is device dependent, we will evaluate the performance of the algorithms for up to 200 per-frame iterations. Fig. 1 shows the average spectral convergence depending on the number of iterations and the number of look-ahead frames $N_{\text{LA}} \in \{0, 1, 2, 3, 7\}$ given in the brackets. $N_{\text{LA}} = 7$ is the maximum number of look-ahead frames which directly overlap with the currently processed frame. The results obtained by the non-iterative RTPGHI algorithm are depicted as horizontal dashed lines. Note that we are intentionally using a fixed range of values

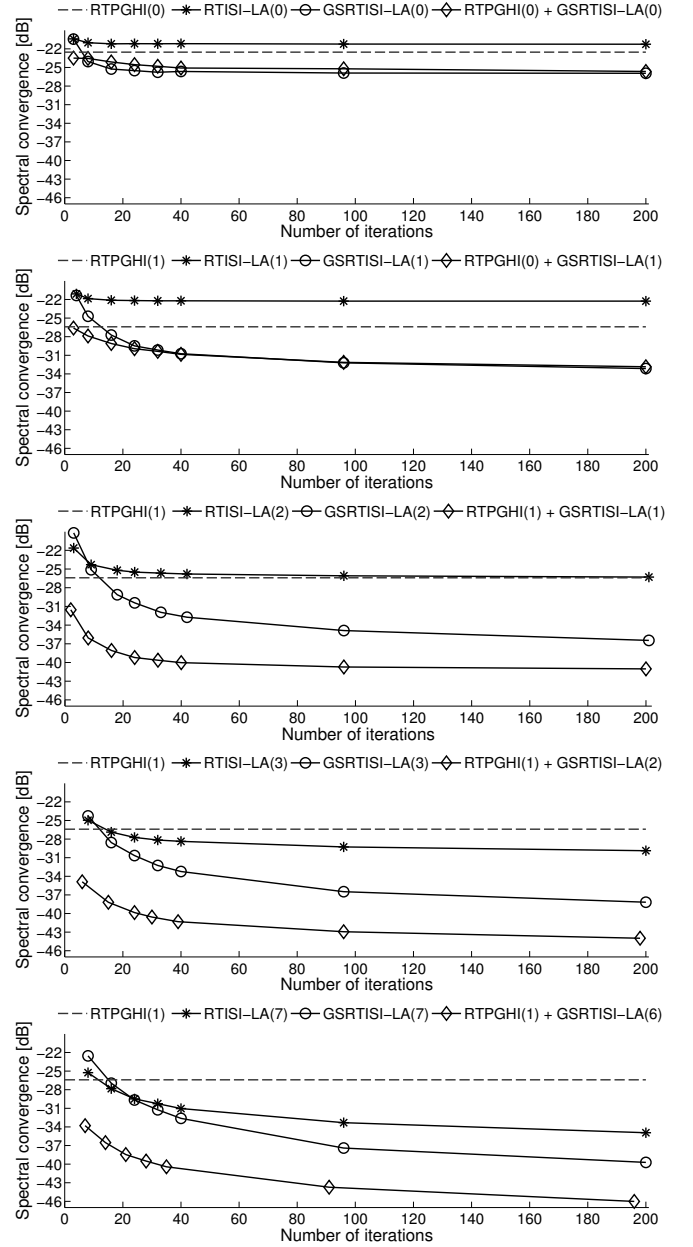


Fig. 1: Comparison of algorithms.

on the y axes. The setting was chosen as the one used in [25, Fig. 6a] to allow a direct comparison. Note that the proposed algorithm with $N_{\text{LA}} = 7$ (corresponds to 87 ms input-output delay) outperforms even the best of the offline-only algorithms.

One can observe that a common behavior of the iterative algorithms is that the average spectral convergence initially decreases rapidly and from a certain number of iterations it starts to level off. This phenomenon could be explained by the fact that some signals in the database reach “convergence” at some point while others continue to improve. Further, one can observe that the proposed algorithm clearly outperforms others in settings using 2 or more look-ahead frames. The scores for individual files for $N_{\text{LA}} = 2$ ($N_{\text{LA}} = 1$ in case of RTPGHI) and 24 per-frame iterations, as well as sound examples for all the samples from the SQAM database, can

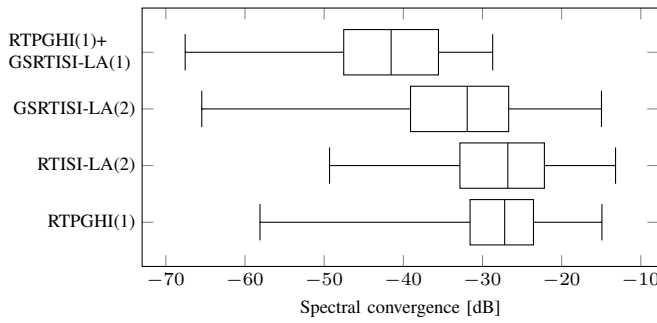


Fig. 2: Box plot of the errors obtained for $N_{LA} = 2$ ($N_{LA} = 1$ in case of RTPGHI). The whiskers denote the minimum and the maximum from the 70 sound excerpts.

be found at the accompanying web page <http://lftat.github.io/notes/048>. Additionally, a box plot of the results is depicted in Fig. 2.

When inspecting results obtained for individual sound excerpts, one can notice that the iterative algorithms struggle with reconstructing recordings of percussion instruments such as claves and castanets and with attacks of transients in general. Conveniently, the RTPGHI algorithm performs very well in such cases and the combination with GSRTISI-LA inherits and even improves upon the behavior as indicated by the low maximum error in Fig. 2.

A real-time demo allowing one-to-one comparison of the algorithms is available in the PHASERET toolbox as `demo_blockproc_phaseret2.m`.

V. CONCLUSION

It has been shown that the combination of GSRTISI-LA and RTPGHI outperforms either of the individual algorithms as well as the RTISI-LA algorithm as soon as enough look-ahead frames is used (see Fig. 1).

Although we have only shown objective error measures in this letter, according to our experience, the quality of the reconstructed signal reflects the error measure improvement. An interested reader can verify this claim by listening to the sound samples found at the accompanying webpage or by running `demo_blockproc_phaseret2.m` using his/her custom audio examples.

In this letter we have assumed that the phase is unknown completely and only the original clean magnitude is known. The proposed algorithm can be easily modified to respect and use coefficients with known phase, but, in the real-world, noisy or modified magnitude and phase are usually observed. Therefore, as the future work, we will focus on simultaneous magnitude and phase estimation given corrupted, noisy or incomplete information as the phase-aware signal processing is currently an active topic of research [28], [29], [30].

ACKNOWLEDGMENT

The author thanks Pavel Rajmic, Thibaud Necciari and Nicki Holighaus for their valuable comments.

REFERENCES

- [1] T. Virtanen, "Monaural sound source separation by nonnegative matrix factorization with temporal continuity and sparseness criteria," *Trans. Audio, Speech and Lang. Proc.*, vol. 15, no. 3, pp. 1066–1074, Mar. 2007. [Online]. Available: <http://dx.doi.org/10.1109/TASL.2006.885253>
- [2] J. Laroche and M. Dolson, "Phase-vocoder: about this phasiness business," in *Applications of Signal Processing to Audio and Acoustics, 1997. 1997 IEEE ASSP Workshop on*, Oct 1997, pp. 4 pp.–.
- [3] P. Smaragdis, B. Raj, and M. Shashanka, "Missing data imputation for time-frequency representations of audio signals," *Journal of Signal Processing Systems*, vol. 65, no. 3, pp. 361–370, 2011.
- [4] D. Griffin and J. Lim, "Signal estimation from modified short-time Fourier transform," *Acoustics, Speech and Signal Processing, IEEE Transactions on*, vol. 32, no. 2, pp. 236–243, Apr 1984.
- [5] N. Perraudin, P. Balazs, and P. S ndergaard, "A fast Griffin-Lim algorithm," in *Applications of Signal Processing to Audio and Acoustics (WASPAA), IEEE Workshop on*, Oct 2013, pp. 1–4.
- [6] R. Decorsiere, P. S ndergaard, E. MacDonald, and T. Dau, "Inversion of auditory spectrograms, traditional spectrograms, and other envelope representations," *Audio, Speech, and Language Processing, IEEE/ACM Transactions on*, vol. 23, no. 1, pp. 46–56, Jan 2015.
- [7] J. Le Roux, H. Kameoka, N. Ono, and S. Sagayama, "Fast signal reconstruction from magnitude STFT spectrogram based on spectrogram consistency," in *Proc. 13th Int. Conf. on Digital Audio Effects (DAFx-10)*, Sep. 2010, pp. 397–403.
- [8] X. Zhu, G. T. Beauregard, and L. Wyse, "Real-time signal estimation from modified short-time Fourier transform magnitude spectra," *Audio, Speech, and Language Processing, IEEE Transactions on*, vol. 15, no. 5, pp. 1645–1653, July 2007.
- [9] G. T. Beauregard, M. Harish, and L. Wyse, "Single pass spectrogram inversion," in *Digital Signal Processing (DSP), IEEE International Conference on*, July 2015, pp. 427–431.
- [10] P. Margon, R. Badeau, and B. David, "Phase reconstruction of spectrograms with linear unwrapping: application to audio signal restoration," in *Proc. 23rd European Signal Processing Conference (EUSIPCO 2015)*, Aug 2015.
- [11] M. Chami, J. Di Martino, L. Pierron, and E. H. Ibn Elhaj, "Real-Time Signal Reconstruction from Short-Time Fourier Transform Magnitude Spectra Using FPGAs," in *5th. International Conference on Information Systems and Economic Intelligence - SIIE 2012*, Djerba, Tunisia, Feb. 2012. [Online]. Available: <https://hal.inria.fr/hal-00761783>
- [12] G. T. Beauregard, X. Zhu, and L. Wyse, "An efficient algorithm for real-time spectrogram inversion," in *Proc. 8th International Conference on Digital Audio Effects (DAFx-05)*, Sep. 2005.
- [13] X. Zhu, G. T. Beauregard, and L. Wyse, "Real-time iterative spectrum inversion with look-ahead," in *Proc. IEEE International Conference on Multimedia and Expo*, 2006.
- [14] V. Gnann and M. Spiertz, "Comb-filter free audio mixing using STFT magnitude spectra and phase estimation," in *Proc. 11th Int. Conf. on Digital Audio Effects (DAFx-08)*, Sep. 2008.
- [15] —, "Inversion of STFT magnitude spectrograms with adaptive window lengths," in *Proc. of IEEE International Conference on Acoustics, Speech and Signal Processing ICASSP '09*, Apr. 2009, pp. 325–328.
- [16] —, "Improving RTISI phase estimation with energy order and phase unwrapping," in *Proc. 13th International Conference on Digital Audio Effects (DAFx-10)*, Sep. 2010.
- [17] Z. Pr  sa and P. L. S ndergaard, "Real-Time Spectrogram Inversion Using Phase Gradient Heap Integration," in *Proc. Int. Conf. Digital Audio Effects (DAFx-16)*, Sep 2016.
- [18] J. W. Eaton, D. Bateman, S. Hauberg, and R. Wehbring, *GNU Octave version 4.0.0 manual: A high-level interactive language for numerical computations*, 2015. [Online]. Available: <http://www.gnu.org/software/octave/doc/interpreter>
- [19] P. L. S ndergaard, B. Torr  sani, and P. Balazs, "The Linear Time Frequency Analysis Toolbox," *International Journal of Wavelets, Multiresolution Analysis and Information Processing*, vol. 10, no. 4, 2012.
- [20] Z. Pr  sa, P. L. S ndergaard, N. Holighaus, C. Wiesmeyer, and P. Balazs, "The Large Time-Frequency Analysis Toolbox 2.0," in *Sound, Music, and Motion*, ser. Lecture Notes in Computer Science. Springer International Publishing, 2014, pp. 419–442.
- [21] T. Strohmer, *Numerical Algorithms for Discrete Gabor Expansions*. Birkh  user Boston, 1998, ch. 8, pp. 267–294.
- [22] P. L. S ndergaard, "Finite discrete Gabor analysis," Ph.D. dissertation, Technical University of Denmark, 2007, available from: <http://lftat.github.io/notes/lftatnote003.pdf>.

- [23] ———, “Efficient algorithms for the discrete Gabor transform with a long FIR window,” *J. Fourier Anal. Appl.*, vol. 18, no. 3, pp. 456–470, 2012.
- [24] I. Daubechies, A. Grossmann, and Y. Meyer, “Painless nonorthogonal expansions,” *Journal of Mathematical Physics*, vol. 27, no. 5, pp. 1271–1283, 1986.
- [25] Z. Průša, P. Balazs, and P. L. Søndergaard, “A Non-iterative Method for STFT Phase Reconstruction,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2016, in preparation. Preprint will be available at <http://ltfat.github.io/notes/ltfatnote040.pdf>.
- [26] N. Sturmel and L. Daudet, “Signal reconstruction from STFT magnitude: A state of the art,” *Proc. 14th Int. Conf. Digital Audio Effects (DAFx-11)*, pp. 375–386, 2011.
- [27] “Tech 3253: Sound Quality Assessment Material recordings for subjective tests,” The European Broadcasting Union, Geneva, Tech. Rep., Sept. 2008. [Online]. Available: <https://tech.ebu.ch/docs/tech/tech3253.pdf>
- [28] T. Gerkmann, M. Krawczyk-Becker, and J. Le Roux, “Phase processing for single-channel speech enhancement: History and recent advances,” *Signal Processing Magazine, IEEE*, vol. 32, no. 2, pp. 55–66, March 2015.
- [29] P. Mowlaee, J. Kulmer, J. Stahl, and F. Mayer, *Single Channel Phase-Aware Signal Processing in Speech Communication: Theory and Practice*. John Wiley & Sons, Inc., 2016.
- [30] P. Mowlaee, R. Saeidi, and Y. Stylianou, “Advances in phase-aware signal processing in speech communication,” *Speech Communication*, vol. 81, pp. 1 – 29, 2016.