

スペクトログラム無矛盾性を用いた独立低ランク行列分析*

☆豊島直, 北村大地 (香川高専), 矢田部浩平 (早稲田大)

1 はじめに

音源分離とは、複数の音源が混合した観測信号から、混合前の音源を推定する技術である。特に、音源位置やマイクロホン位置等が未知の条件で音源分離を達成する技術はブラインド音源分離 (blind source separation: BSS) と呼ばれる。観測信号のチャンネル数 (マイクロホン数) と混合されている音源数が等しい条件下では、独立成分分析 (independent component analysis: ICA) [1] に基づく BSS として、周波数領域 ICA (frequency-domain ICA: FDICA) [2], 独立ベクトル分析 (independent vector analysis: IVA) [3, 4], 及び独立低ランク行列分析 (independent low-rank matrix analysis: ILRMA) [5, 6] 等が提案されてきた。これらはいずれも、FDICA におけるパーミュテーション問題 [7] (周波数毎に得られる分離信号の順序を整列する問題) を解決するための FDICA の拡張である。IVA や ILRMA は、その分離性能の高さや初期値に対する頑健性等の利点から、数多くの改良手法や一般化手法へと発展している (例えば [8] 等)。

一方、多くの音響信号処理で用いられる短時間フーリエ変換 (short-time Fourier transform: STFT) では、窓関数とそのオーバーラップにより近傍時間周波数グリッドに共起性が生じる (文献 [9] の Fig. 1 参照)。これは時間周波数表現された信号が根本的に有している性質であり、スペクトログラム無矛盾性と呼ばれる [10, 11]。近年、このスペクトログラム無矛盾性を FDICA や IVA に導入した BSS が提案され、分離行列最適化時のパーミュテーション問題の解決をアシストする新しい基準となることが示された [9]。その原理は、分離信号のパーミュテーションが近傍周波数で異なる場合に、スペクトログラム無矛盾性 (隣接周波数の共起関係等) が大きく損なわれるという性質に基づいている (文献 [9] の Fig. 2 参照)。

本稿では、文献 [9] で示された結果に基づき、より高精度な BSS 達成を目的として、スペクトログラム無矛盾性を ILRMA に導入したアルゴリズムを提案する。また、音楽信号と音声信号の BSS 実験を通して、ILRMA においてもスペクトログラム無矛盾性が分離性能の向上に貢献することを示す。

2 BSS とスペクトログラム無矛盾性

2.1 STFT 及び BSS の定式化

離散時間信号の l 番目のサンプルを $x[l]$ のように表記し、 N 個の音源信号が M 個のマイクロホンで観測される状況を考える。多チャンネルの音源信号、観測信号、及び分離信号をそれぞれ次式で表す。

$$\mathbf{s}[l] = [s_1[l], \dots, s_n[l], \dots, s_N[l]]^T \in \mathbb{R}^N \quad (1)$$

$$\mathbf{x}[l] = [x_1[l], \dots, x_m[l], \dots, x_M[l]]^T \in \mathbb{R}^M \quad (2)$$

$$\mathbf{y}[l] = [y_1[l], \dots, y_n[l], \dots, y_N[l]]^T \in \mathbb{R}^N \quad (3)$$

ここで、 $n = 1, 2, \dots, N$, $m = 1, 2, \dots, M$, 及び $l = 1, 2, \dots, L$ はそれぞれ音源、マイクロホン (チャ

ネル)、及び離散時間のインデクスであり、 \cdot^T は転置を表す。BSS では、音源信号 \mathbf{s} に近い分離信号 \mathbf{y} を、観測信号 \mathbf{x} から推定することが目的となる。

FDICA, IVA, 及び ILRMA 等の周波数領域 BSS では、信号を時間周波数領域で取り扱う。STFT における窓長とソフト長をそれぞれ Q 及び τ とおくと、信号 $z[l]$ の j 番目の短時間セグメントは次式となる。

$$\begin{aligned} z^{[j]} &= [z[(j-1)\tau+1], z[(j-1)\tau+2], \\ &\quad \dots, z[(j-1)\tau+Q]]^T \\ &= [z^{[j]}[1], \dots, z^{[j]}[q], \dots, z^{[j]}[Q]]^T \in \mathbb{R}^Q \end{aligned}$$

ここで、 $j = 1, 2, \dots, J$ 及び $q = 1, 2, \dots, Q$ はそれぞれ短時間セグメント及びセグメント内サンプルのインデクスを表す。セグメント数は (必要に応じて信号に零詰めを行ったうえで) $J = L/\tau$ となる。信号 $\mathbf{z} = [z[1], \dots, z[l], \dots, z[L]]^T \in \mathbb{R}^L$ の STFT を

$$\mathbf{Z} = \text{STFT}_\omega(\mathbf{z}) \in \mathbb{C}^{I \times J} \quad (4)$$

と表し、 \mathbf{Z} の (i, j) 要素は次式で与えられる。

$$z_{ij} = \sum_{q=1}^Q \omega[q] z^{[j]}[q] e^{-i2\pi(q-1)(i-1)/F} \quad (5)$$

ここで、 $i = 1, 2, \dots, I$ は周波数ピンのインデクス、 F は $\lfloor F/2 \rfloor + 1 = I$ を満たす整数、 $\lfloor \cdot \rfloor$ は床関数、 i は虚数単位、及び ω は解析時の窓関数を表す。合成時の窓関数を $\tilde{\omega}$ とおくと、逆 STFT を $\text{ISTFT}_{\tilde{\omega}}(\cdot)$ と表記する。本稿では、 ω と $\tilde{\omega}$ のペアが次式の完全再構成条件を満たすことを仮定する。

$$\mathbf{z} = \text{ISTFT}_{\tilde{\omega}}(\text{STFT}_\omega(\mathbf{z})) \quad \forall \mathbf{z} \in \mathbb{R}^L \quad (6)$$

各チャンネルに STFT を適用して得られる音源信号、観測信号、及び分離信号のスペクトログラムの (i, j) 番目の要素をそれぞれ次式で表す。

$$\mathbf{s}_{ij} = [s_{ij1}, \dots, s_{ijn}, \dots, s_{ijN}]^T \in \mathbb{C}^N \quad (7)$$

$$\mathbf{x}_{ij} = [x_{ij1}, \dots, x_{ijm}, \dots, x_{ijM}]^T \in \mathbb{C}^M \quad (8)$$

$$\mathbf{y}_{ij} = [y_{ij1}, \dots, y_{ijn}, \dots, y_{ijN}]^T \in \mathbb{C}^N \quad (9)$$

また、式 (7)–(9) の時間周波数行列をそれぞれ $\mathbf{S}_n \in \mathbb{C}^{I \times J}$, $\mathbf{X}_m \in \mathbb{C}^{I \times J}$, 及び $\mathbf{Y}_n \in \mathbb{C}^{I \times J}$ と定義する。周波数領域 BSS では、次式の瞬時混合を仮定する。

$$\mathbf{x}_{ij} = \mathbf{A}_i \mathbf{s}_{ij} \quad (10)$$

ここで、 $\mathbf{A}_i \in \mathbb{C}^{M \times N}$ は周波数毎の混合行列である。式 (10) は、混合系の残響時間が窓長よりも十分短い場合に近似的に成立する。

以後、決定的な系 ($M = N$) を考える。この場合、BSS は \mathbf{A}_i の逆行列を推定する問題となる。この逆行列を $\mathbf{W}_i \approx \mathbf{A}_i^{-1}$ とすると、分離信号は次式となる。

$$\mathbf{y}_{ij} = \mathbf{W}_i \mathbf{x}_{ij} \quad (11)$$

*Independent low-rank matrix analysis based on spectrogram consistency. By Nao TOSHIMA, Daichi KITAMURA (NIT Kagawa), and Kohei YATABE (Waseda Univ.).

ここで、 $\mathbf{W}_i = [\mathbf{w}_{i1}, \dots, \mathbf{w}_{in}, \dots, \mathbf{w}_{iN}]^H \in \mathbb{C}^{N \times M}$ は分離行列と呼ばれ、 \cdot^H はエルミート転置を表す。

ICAに基づくBSSでは、分離音源のスケールや順序が不定であるため、次式の任意性が存在する。

$$\hat{y}_{ij} = \hat{\mathbf{W}}_i \mathbf{x}_{ij} \quad (\hat{\mathbf{W}}_i = \mathbf{D}_i \mathbf{P}_i \mathbf{W}_i) \quad (12)$$

ここで、 $\mathbf{D}_i \in \mathbb{C}^{N \times N}$ 及び $\mathbf{P}_i \in \{0, 1\}^{N \times N}$ はそれぞれ任意の対角行列及びパーミュテーション行列である。周波数毎のスケールの任意性についてはプロジェクトンバック法 [13] によって解析的に復元可能であるが、周波数毎の分離信号の順序の整列はパーミュテーション問題と呼ばれ、大きな課題である。

2.2 スペクトログラム無矛盾性に基づくパーミュテーション解決 [9]

文献 [9] では、スペクトログラム無矛盾性と呼ばれる STFT の性質を用いてパーミュテーション問題が解決できる可能性が示されている。スペクトログラムは本来、時間周波数表現における不確定性原理より、近傍の時間周波数成分が共起関係を持つ。スペクトログラム \mathbf{Z} において、この共起関係が崩れている状態を矛盾と呼び、逆 STFT はその矛盾した成分を無矛盾な状態に復元する。即ち、 \mathbf{Z} の無矛盾性は

$$\mathcal{E}(\mathbf{Z}) = \mathbf{Z} - \text{STFT}_\omega(\text{ISTFT}_\omega(\mathbf{Z})) \quad (13)$$

のノルム $\|\mathcal{E}(\mathbf{Z})\|$ によって特徴付けられ、それが 0 となるスペクトログラムを無矛盾と呼ぶ。大雑把に言えば、逆 STFT はスペクトログラムの時間及び周波数方向への滲みがない成分に対して、スムージングをかけるような処理となる（文献 [9] の Fig. 1 参照）。

文献 [9] の Fig. 2 のように、パーミュテーション問題が生じている分離信号のスペクトログラムは、隣接周波数の不連続性に起因して無矛盾性が大きく損なわれている。従って、FDICA や IVA における \mathbf{W}_i の反復最適化の過程で $\text{STFT}_\omega(\text{ISTFT}_\omega(\mathbf{Y}_n))$ を作用させることで、周波数毎の音源成分が正しく整列された（パーミュテーション問題が解決された）状態に少しずつ誘導できることが示されている。本稿では、より高い分離性能を発揮できる ILRMA においても、上記のスペクトログラム無矛盾性を反復毎に担保するアルゴリズムが性能向上に貢献することを実証する。

3 提案手法

3.1 スペクトログラム無矛盾性に基づく ILRMA

ILRMA は次式の目的関数を最小化する [6]。

$$\mathcal{L} = -2J \sum_i |\det \mathbf{W}_i| + \sum_{i,j,n} \left(\frac{|\mathbf{w}_{in}^H \mathbf{x}_{ij}|^2}{\sum_k t_{ikn} v_{kjn}} + \log \sum_k t_{ikn} v_{kjn} \right) \quad (14)$$

ここで、 t_{ikn} 及び v_{kjn} は非負行列 $\mathbf{T}_n \in \mathbb{R}_{\geq 0}^{I \times K}$ 及び $\mathbf{V}_n \in \mathbb{R}_{\geq 0}^{K \times J}$ の要素であり、 $k = 1, 2, \dots, K$ は \mathbf{T}_n の列ベクトルのインデクスである。式 (14) の最小化は、ICA に基づく分離信号間の独立性の最大化と非負値行列因子分解 (nonnegative matrix factorization: NMF) [12] に基づく分離信号の低ランク近似 $|\mathbf{Y}_n|^2 \approx \mathbf{T}_n \mathbf{V}_n$ を同時に行うことに対応する。但し、行列に対する $|\cdot|^2$ は要素毎の絶対値二乗を表す。これは即ち、

混合前の音源信号 \mathbf{S}_n が低ランク時間周波数構造を有することを仮定し、分離信号 \mathbf{Y}_n を NMF でモデル化しながら分離行列 \mathbf{W}_i を最適化することで、パーミュテーション問題が生じない \mathbf{W}_i を推定している。

提案手法では、スペクトログラム無矛盾性を担保する処理を ILRMA の最適化アルゴリズムに組み込む。従来の ILRMA の分離行列の反復最適化計算 [4]

$$\mathbf{U}_{in} \leftarrow \frac{1}{J} \sum_j \frac{1}{\sum_k t_{ikn} v_{kjn}} \mathbf{x}_{ij} \mathbf{x}_{ij}^H \quad (15)$$

$$\mathbf{w}_{in} \leftarrow (\mathbf{W}_i \mathbf{U}_{in})^{-1} \mathbf{e}_n \quad (16)$$

$$\mathbf{w}_{in} \leftarrow \mathbf{w}_{in} (\mathbf{w}_{in}^H \mathbf{U}_{in} \mathbf{w}_{in})^{-\frac{1}{2}} \quad (17)$$

$$y_{ijn} \leftarrow \mathbf{w}_{in}^H \mathbf{x}_{ij} \quad (18)$$

及び NMF 音源モデルの反復最適化計算

$$t_{ikn} \leftarrow t_{ikn} \sqrt{\frac{\sum_j |y_{ijn}|^2 (\sum_{k'} t_{ik'n} v_{k'jn})^{-2} v_{kjn}}{\sum_j (\sum_{k'} t_{ik'n} v_{k'jn})^{-1} v_{kjn}}} \quad (19)$$

$$v_{kjn} \leftarrow v_{kjn} \sqrt{\frac{\sum_i |y_{ijn}|^2 (\sum_{k'} t_{ik'n} v_{k'jn})^{-2} t_{ikn}}{\sum_i (\sum_{k'} t_{ik'n} v_{k'jn})^{-1} t_{ikn}}} \quad (20)$$

において、

$$\mathbf{Y}_n \leftarrow \text{STFT}_\omega(\text{ISTFT}_\omega(\mathbf{Y}_n)) \quad (21)$$

なる演算を挿入することで、毎回の反復においてスペクトログラム無矛盾性を担保する。ここで、 $\mathbf{e}_n \in \{0, 1\}^N$ は n 番目の要素のみが 1 の単位ベクトルである。式 (21) は、分離信号のスペクトログラム \mathbf{Y}_n を無矛盾スペクトログラムの集合へと射影していることに対応する。そのため、もし \mathbf{Y}_n が無矛盾であれば式 (21) は何もしておらず、 \mathbf{Y}_n に矛盾があれば時間及び周波数の両方向にスムージングがかかる。

3.2 反復毎のプロジェクトンバック法の適用

ILRMA の推定は FDICA や IVA と同様に、式 (12) に示す任意性がある。とくに、任意の対角行列 \mathbf{D}_i に起因する周波数毎のスケール不定性は、それ自身がスペクトログラム無矛盾性を崩してしまう。提案手法では、 \mathbf{D}_i に起因する矛盾成分の影響を最小限に抑えるため、毎回の反復計算において式 (21) を行う直前に、次式のプロジェクトンバック法 [13] を適用する。

$$\tilde{y}_{ijn} = \mathbf{W}_i^{-1} (\mathbf{e}_n \circ \mathbf{y}_{ij}) = y_{ijn} \lambda_{in}, \quad (22)$$

ここで、 $\tilde{\mathbf{y}}_{ijn} = [\tilde{y}_{ijn1}, \dots, \tilde{y}_{ijnm}, \dots, \tilde{y}_{ijnM}]^T \in \mathbb{C}^M$ はスケール補正後の分離信号の (i, j) 番目の成分、 $\lambda_{in} = [\lambda_{in1}, \dots, \lambda_{inm}, \dots, \lambda_{inM}]^T \in \mathbb{C}^M$ はスケール補正係数、及び \circ は要素毎の積を表す。また、式 (22) によって目的関数 (14) の値が変動することを防ぐために、他の変数も次のように補正する。

$$\mathbf{w}_{in} \leftarrow \mathbf{w}_{in} \lambda_{inm_{\text{ref}}} \quad (23)$$

$$y_{ijn} \leftarrow \mathbf{w}_{in}^H \mathbf{x}_{ij} \quad (24)$$

$$t_{ikn} \leftarrow t_{ikn} |\lambda_{inm_{\text{ref}}}|^2 \quad (25)$$

ここで、 m_{ref} はプロジェクトンバック法で用いるリファレンスチャンネルのインデクスである。

これらをまとめると、提案手法は Algorithm 1 に示すアルゴリズムによって実現される。

Algorithm 1 Consistent ILRMA

Input: $\{\mathbf{x}_{ij}\}_{i=1,j=1}^{I,J}$, maxIter

Output: $\{\mathbf{y}_{ij}\}_{i=1,j=1}^{I,J}$

- 1: Initialize $\{\mathbf{T}_n\}_{n=1}^N, \{\mathbf{V}_n\}_{n=1}^N, \{\mathbf{W}_i\}_{i=1}^I$
- 2: **for** iter = 1, 2, \dots , maxIter **do**
- 3: Ensure consistency by (21) $\forall n$
- 4: Update \mathbf{W}_i by (15)–(18) $\forall i, j, n$
- 5: Update \mathbf{T}_n and \mathbf{V}_n by (19), (20) $\forall i, j, k, n$
- 6: Apply back projection by (22) $\forall i, j, n$
- 7: Update parameters by (23)–(25) $\forall i, j, k, n$
- 8: **end for**

Table 1 Experimental conditions

Window function	Hann window
Window length	64, 128, 256, 512, 768, 1024 ms
Window shift length	1/16 or 1/2 of window length
Number of bases K for each source in ILRMA	10 for music signals and 2 for speech signals
Initialization of parameters	\mathbf{W}_i : identity matrix \mathbf{T}_n and \mathbf{V}_n : random matrices
Number of iterations	100
Number of trials	5 with different random seeds
Reference channel m_{ref}	1

4 BSS 実験

4.1 実験条件

提案手法の有効性を確認するために、BSS性能を従来のILRMAと比較した。ここで、従来手法は提案手法から式(21)の演算のみを省いたものとし、その他の条件は全て提案手法と統一している。実験には、文献[14]と同じ音楽信号4種と音声信号4種及び混合系(インパルス応答)を用いた。本稿ではインパルス応答JR2($T_{60} = 470$ ms)の結果のみを示すが、文献[15]には網羅的な実験結果が示されているので参照されたい。評価値は音源対歪み比(source-to-distortion ratio: SDR)[16]の改善量を用いた。その他の実験条件はTable 1に示す通りである。

4.2 実験結果

Fig. 1は、提案手法の各反復における矛盾成分のエネルギー値の例である。ここで、縦軸 $\|\mathcal{E}(\mathbf{Y})\|_2^2/\|\mathbf{X}\|_2^2$ 中の \mathbf{X} 及び \mathbf{Y} はそれぞれ $\mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2]$ 及び $\mathbf{Y} = [\mathbf{Y}_1, \mathbf{Y}_2]$ を表し、 \mathcal{E} の定義は式(13)に示されている。矛盾成分のエネルギー値は反復初期において大きく増加するものの、反復後半では一定の値に収束している。このことから、提案手法の最適化ができるだけ無矛盾な解へと誘導する様子を確認できる。

Fig. 2は、各窓長及びシフト長での平均SDR改善量の比較を示している。図中のラベルConv.及びProp.はそれぞれ従来手法と提案手法に対応しており、異なる初期値に対する試行や4種類の信号に対する結果を全てまとめて箱ひげ図で示している。この結果より、式(10)が成立ししやすい長い窓長においてILRMAは高い性能を発揮することが確認できる。また提案手法は、従来手法の分離性能が比較的高い条件において性能が向上する傾向が確認でき、場合によっては2~3 dB程度の改善が見られた。これは、音源分離が成功するほど、分離信号 \mathbf{Y}_n は本来の音源信号 \mathbf{S}_n に近づき、提案手法においてスペクトログラム無矛盾性を担保する効果が高くなるためと推測される。

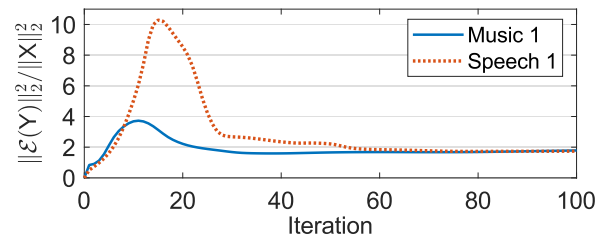


Fig. 1 Normalized energy of the inconsistent components of proposed method, where window and shift lengths were 256 ms and 32 ms, respectively.

5 おわりに

本稿では、反復最適化の中でスペクトログラム無矛盾性を担保するILRMAを提案し、その有効性を実験的に確認した。今後の課題として、スペクトログラム無矛盾性を陽に加味した新しい音源モデルの導入や目的関数の設計等が挙げられる。

謝辞 本研究の一部はJSPS 科研費19K20306の助成を受けたものである。

参考文献

- [1] P. Comon, "Independent component analysis, a new concept?," *Signal Process.*, vol. 36, no. 3, pp. 287–314, 1994.
- [2] P. Smaragdis, "Blind separation of convolved mixtures in the frequency domain," *Neurocomputing*, vol. 22, pp. 21–34, 1998.
- [3] T. Kim, H. T. Attias, S.-Y. Lee, and T.-W. Lee, "Blind source separation exploiting higher-order frequency dependencies," *IEEE Trans. ASLP*, vol. 15, no. 1, pp. 70–79, 2007.
- [4] N. Ono, "Stable and fast update rules for independent vector analysis based on auxiliary function technique," *Proc. WASPAA*, pp. 189–192, 2011.
- [5] 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. ASLP*, vol. 24, no. 9, pp. 1626–1641, 2016.
- [6] D. Kitamura, N. Ono, H. Sawada, H. Kameoka, and H. Saruwatari, "Determined blind source separation with independent low-rank matrix analysis," in *Audio Source Separation*, S. Makino, Ed., pp. 125–155. Springer, Cham, 2018.
- [7] H. Sawada, R. Mukai, S. Araki, and S. Makino, "A robust and precise method for solving the permutation problem of frequency-domain blind source separation," *IEEE Trans. SAP*, vol. 12, no. 5, pp. 530–538, Sep. 2004.
- [8] K. Yatabe and D. Kitamura, "Time-frequency-masking-based determined BSS with application to sparse IVA," *Proc. ICASSP*, pp. 715–719, 2019.
- [9] K. Yatabe, "Consistent ICA: Determined BSS meets spectrogram consistency," *IEEE Signal Process. Lett.*, vol. 27, pp. 870–874, 2020.
- [10] J. Le Roux and E. Vincent, "Consistent Wiener filtering for audio source separation," *IEEE Signal Process. Lett.*, vol. 20, no. 3, pp. 217–220, 2013.
- [11] K. Yatabe, Y. Masuyama, T. Kusano, and Y. Oikawa, "Representation of complex spectrogram via phase conversion," *Acoust. Sci. & Tech.*, vol. 40, no. 3, pp. 170–177, 2019.
- [12] C. Févotte, N. Bertin, and J.-L. Durrieu, "Nonnegative matrix factorization with the Itakura-Saito divergence. With application to music analysis," *Neural Computation*, vol. 21, no. 3, pp. 793–830, 2009.
- [13] K. Matsuoka and S. Nakashima, "Minimal distortion principle for blind source separation," *Proc. ICA*, pp. 722–727, 2001.
- [14] D. Kitamura, N. Ono, and H. Saruwatari, "Experimental analysis of optimal window length for independent low-rank matrix analysis," *Proc. EUSIPCO*, pp. 1210–1214, 2017.
- [15] D. Kitamura and K. Yatabe, "Consistent independent low-rank matrix analysis for determined blind source separation," *arXiv:2007.00274*, 2020.
- [16] E. Vincent, R. Gribonval, and C. Févotte, "Performance measurement in blind audio source separation," *IEEE Trans. ASLP*, vol. 14, no. 4, pp. 1462–1469, 2006.

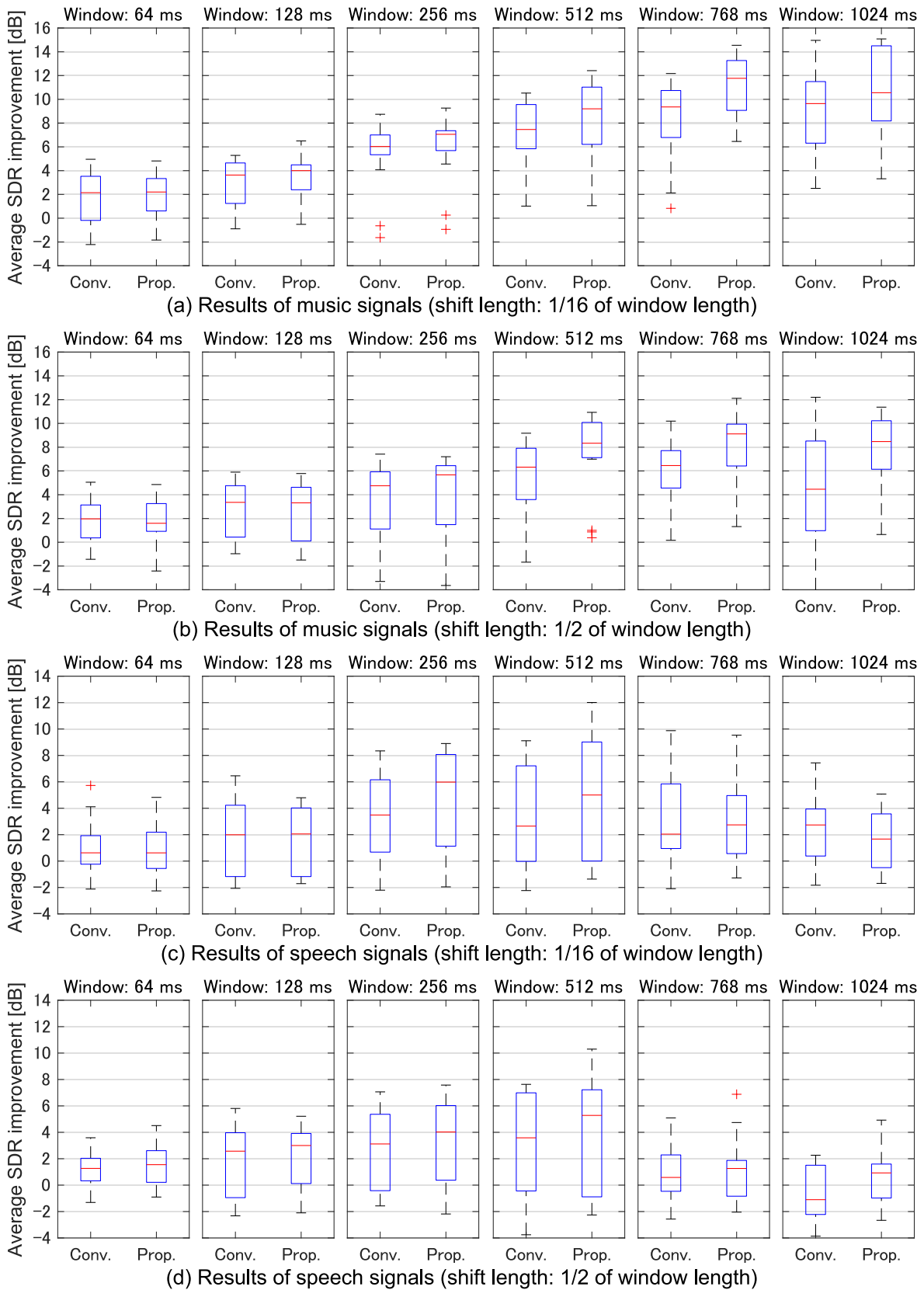


Fig. 2 SDR improvements for (a) music with 1/16 window shifting, (b) music with 1/2 window shifting, (c) speech with 1/16 window shifting, and (d) speech with 1/2 window shifting, where “Conv.” and “Prop.” respectively indicate conventional and proposed methods.