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(香川高専), 北村大地(香川高専), 矢田部浩平(早稲田大)

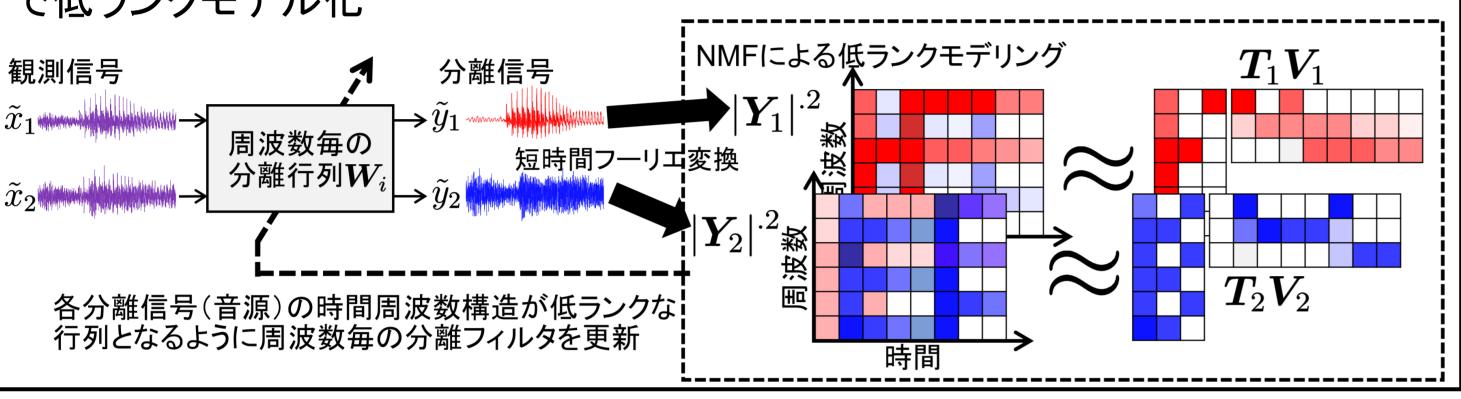


音源分離:複数の音声や楽器音の混合音から個々の音源を推定する技術



応用例

- 自動採譜の前段処理
- 音声認識の精度向上etc.
- ブラインド音源分離(BSS):マイクや音源の位置, 混合系が未知の音源分離 混合信号 分離信号 音源信号 独立成分分析(ICA)[Common, 1994] に基づく手法が一般的 分離音yを求める為の 分離行列Wを推定 混合系 分離系
- 独立低ランク行列分析(independent low-rank matrix analysis: ILRMA):[Kitamura+, 2016] 周波数ビン毎のICA+各音源の時間周波数構造を非負値行列因子分解(NMF) で低ランクモデル化



解決すべき問題

 ILRMA の周波数毎の推定分離行列 $oldsymbol{W}_i$ で得られる時間周波数領域の分 離信号はスペクトログラム無矛盾性が損なわれている

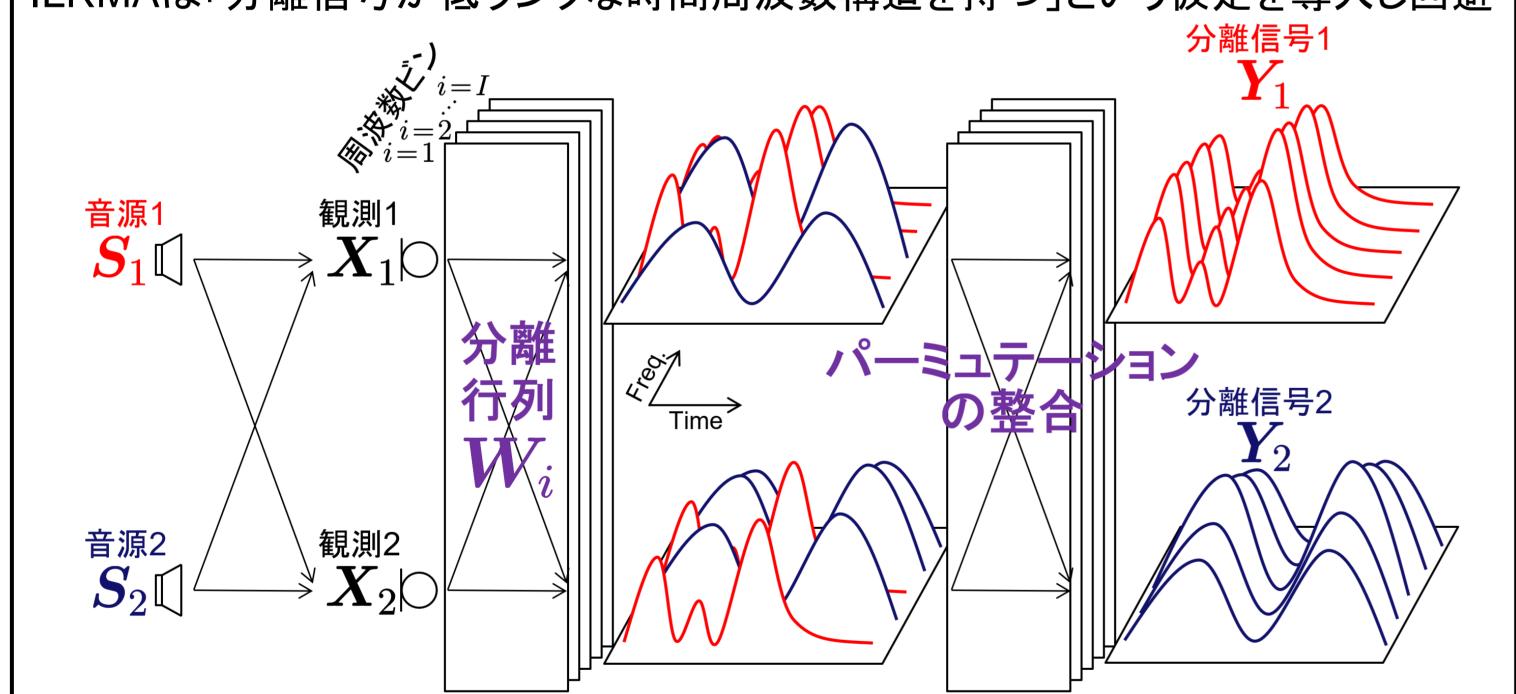
本研究の目的

ILRMAの分離行列推定時にスペクトログラム無矛盾性を担保した新しい アルゴリズムを提案し、分離性能の向上に寄与するかを実験的に調査

2. 従来手法

BSSにおけるパーミュテーション問題

ICAを周波数毎の複素時系列に適用して周波数毎に音源分離することを考える 周波数によって分離信号の音源順序が変わってしまうパーミュテーション問題 ILRMAは「分離信号が低ランクな時間周波数構造を持つ」という仮定を導入し回避



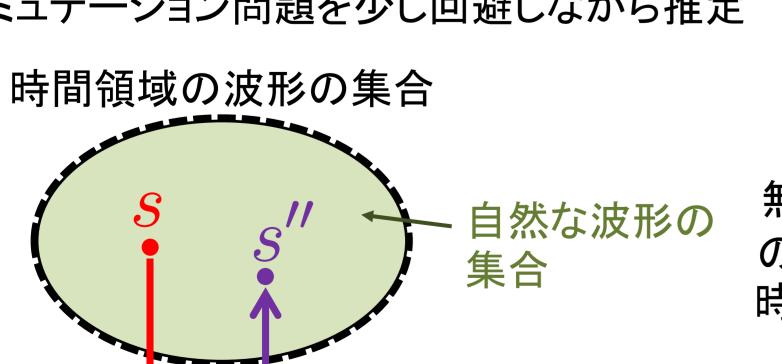
BSSにスペクトログラム無矛盾性を導入 [K.Yatabe, 2020]

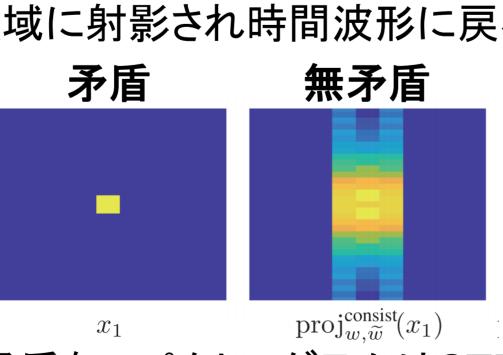
時間波形を短時間フーリエ変換(STFT)するとスペクトログラムが得られ、スペクト ログラムを逆STFT(ISTFT)すると時間波形に戻る(完全再構成条件を満たす場合) スペクトログラムを信号処理で加工すると、矛盾したスペクトログラムとなる 矛盾したスペクトログラムをISTFTすると無矛盾な領域に射影され時間波形に戻る BSSの分離行列の最適化の反復毎に

(時間波形が

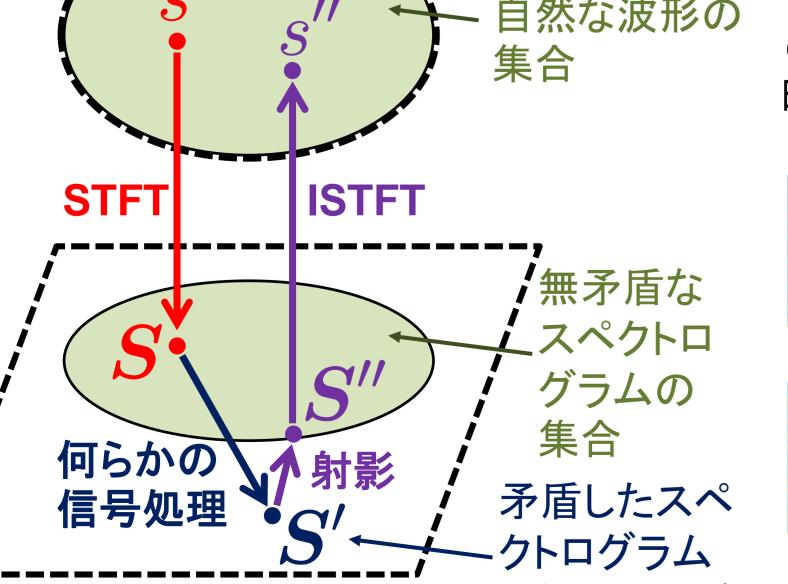
存在しない)

「ISTFT→STFT」という処理を挿入することで パーミュテーション問題を少し回避しながら推定

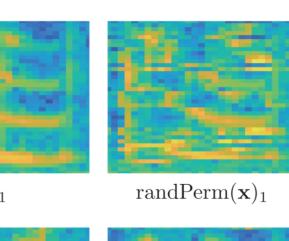


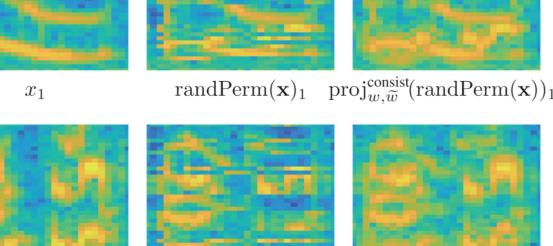


無矛盾なスペクトログラムはSTFT の窓掛けやオーバーラップシフトで 時間周波数の両方向に滲んでいる



時間周波数領域の集合





 $\operatorname{randPerm}(\mathbf{x})_2 \quad \operatorname{proj}_{w.\widetilde{w}}^{\operatorname{consist}}(\operatorname{randPerm}(\mathbf{x}))_2$ パーミュテーション問題を起こすと 矛盾したスペクトログラムになる

3. スペクトログラム無矛盾性に基づくILRMA

ILRMAへのスペクトログラム無矛盾性の適用

Consistent ILRMAのアルゴリズム

 $\mathbf{Input:} \left\{ \boldsymbol{x}_{ij} \right\}_{i=1,j=1}^{I,J}, \mathsf{maxlter}$ Output: $\{oldsymbol{y}_{ij}\}_{i=1,j=1}^{I,J}$ 矛盾した推定分離スペク 1: Initialize $\{T_n\}_{n=1}^N, \{V_n\}_{n=1}^N, \{W_i\}_{i=1}^I$ トログラム Y_n を無矛盾な 2: for iter = $1, 2, \dots, \text{maxIter do}$ 領域へ射影 3: $Y_n \leftarrow \text{STFT}_{\boldsymbol{\omega}}(\text{ISTFT}_{\widetilde{\boldsymbol{\omega}}}(Y_n))$ 原稿のAlgorithm 1に $t_{ikn} \leftarrow t_{ikn} \sqrt{\frac{\sum_{j} |y_{ijn}|^2 \left(\sum_{k'} t_{ik'n} v_{k'jn}\right)^{-2} v_{kjn}}{\sum_{j} \left(\sum_{k'} t_{ik'n} v_{k'jn}\right)^{-1} v_{kjn}}}$ 誤植あり $(15)~(18) \succeq (19)~(20)$ の $v_{kjn} \leftarrow v_{kjn} \sqrt{\frac{\sum_{i} |y_{ijn}|^2 \left(\sum_{k'} t_{ik'n} v_{k'jn}\right)^{-2} t_{ikn}}{\sum_{i} \left(\sum_{k'} t_{ik'n} v_{k'jn}\right)^{-1} t_{ikn}}}$ 順番が逆になっています (20)6: $U_{in} \leftarrow \frac{1}{J} \sum_{i} \frac{1}{\sum_{k} t_{ikn} v_{kjn}} \boldsymbol{x}_{ij} \boldsymbol{x}_{ij}^{\mathrm{H}}$ (15)従来のILRMA 4,5行目でNMF音源 7: $\boldsymbol{w}_{in} \leftarrow (\boldsymbol{W}_{i}\boldsymbol{U}_{in})^{-1} \boldsymbol{e}_{n}$ モデルの最適化 8: $\boldsymbol{w}_{in} \leftarrow \boldsymbol{w}_{in} \left(\boldsymbol{w}_{in}^{\mathrm{H}} \boldsymbol{U}_{in} \boldsymbol{w}_{in} \right)^{-\frac{1}{2}}$ 6~9行目で分離行列 $9: y_{ijn} \leftarrow \boldsymbol{w}_{in}^{\mathrm{H}} \boldsymbol{x}_{ii}$ (18)の最適化

GitHub: https://github.com/d-kitamura/ILRMA/blob/master/consistentILRMA.m

 $m{w}_{in}$ は $m{W}_i$ の行べクトル, y_{ijn} は $m{Y}_n$ の $(i,\ j)$ 要素, t_{ikn} は $m{T}_n$ の $(i,\ k)$ 要素, v_{kjn} は $m{V}_n$ の(k,j)要素, $\lambda_{inm_{\mathrm{ref}}}$ は m_{ref} にプロジェクションバックして求まる補正係数 iは周波数ビンインデクス, jは時間フレームインデクス, kは基底ベクトルインデクス

実験条件

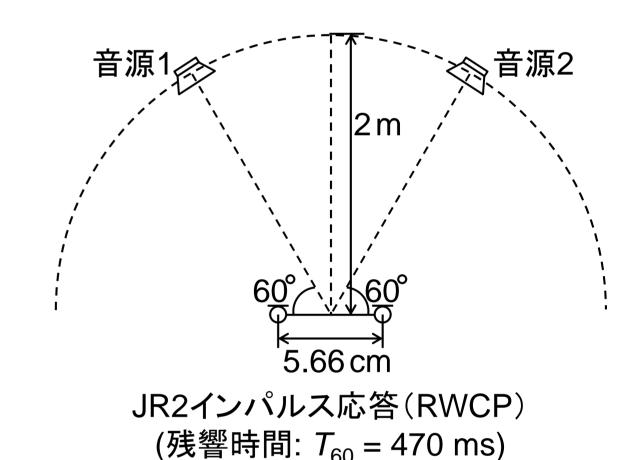
13: **end for**

10: $\boldsymbol{w}_{in} \leftarrow \boldsymbol{w}_{in} \lambda_{inm_{\text{ref}}}$

 $y_{ijn} \leftarrow oldsymbol{w}_{in}^{ ext{H}} oldsymbol{x}_{ij}$

12: $t_{ikn} \leftarrow t_{ikn} |\lambda_{inm_{ref}}|^2$

2 4.55 45 14 11	
窓関数	ハン窓
窓長	64, 128, 256, 512, 768, 1024 ms
シフト長	窓長の1/16 または 1/2
基底数	音楽10, スピーチ2
初期値	$oldsymbol{W}_i$:単位行列
	$oldsymbol{T}_n$ and $oldsymbol{V}_n$: 乱数行列
反復回数	100 回
試行回数	乱数シードを変えて5回
参照マイクチャネル	$m_{\rm ref} = 1$



毎反復においてプロジェ

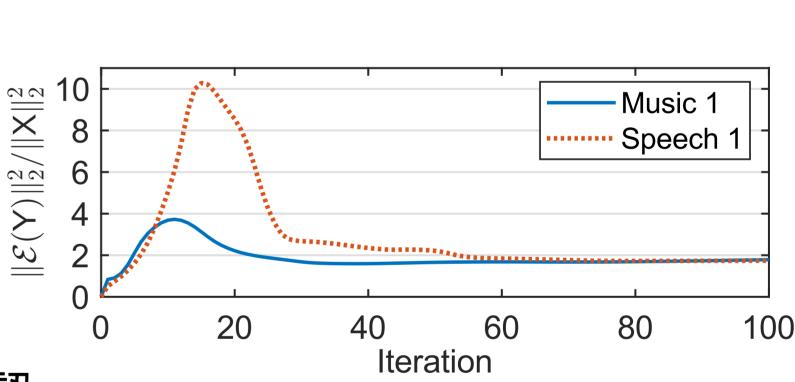
クションバックを適用

実験結果

右図は観測パワーで正規化された 反復毎の矛盾成分量(二乗L2ノルム) $\mathcal{E}(\mathsf{Y}) = \mathsf{Y} - \mathrm{STFT}_{\boldsymbol{\omega}}(\mathrm{ISTFT}_{\widetilde{\boldsymbol{\omega}}}(\mathsf{Y}))$

 $\mathsf{X} = [oldsymbol{X}_1, oldsymbol{X}_2], \ \ \mathsf{Y} = [oldsymbol{Y}_1, oldsymbol{Y}_2]$

一度大きな値に上昇するが、最適化の 後半では一定の値に収束することを確認



右下図は従来手法と提案手法(前章のアルゴリズムの3行目の有無)の比較

音楽信号の考察

STFTの窓長が長い場合 提案手法が明らかに 従来手法を上回る

シフト長の違いによる 傾向差はあまり確認 できず

音声信号の考察

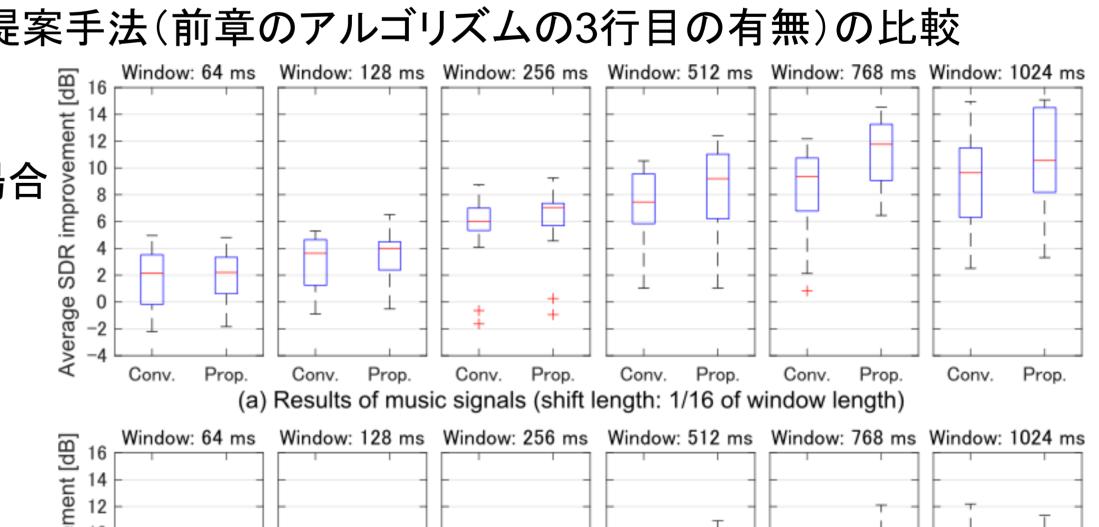
STFTの窓長が512 ms の場合に提案手法が 従来手法を上回る

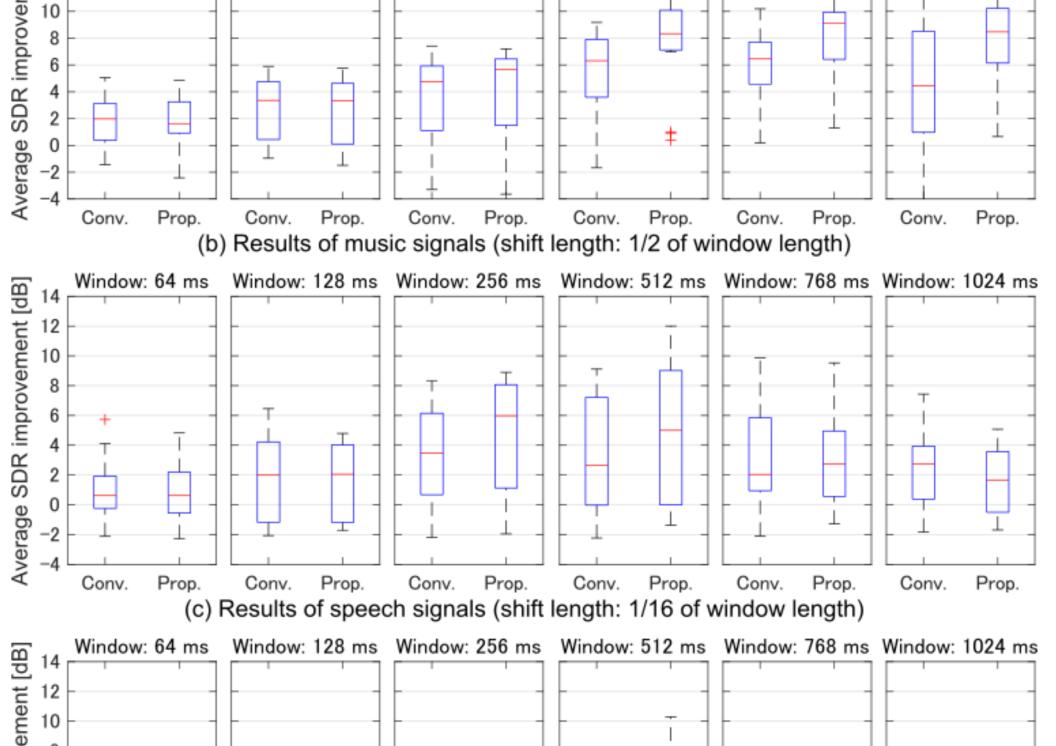
シフト長の違いによる 傾向差はあまり確認 できず

全体的な考察

音源分離が成功する程 提案手法の有効性が 顕著になることを確認

分離が成功した場合, 推定スペクトログラム Y_n は無矛盾なスペクトロ グラムに近づくためと 推測される





(d) Results of speech signals (shift length: 1/2 of window length)