

# 非ガウス性モノラル音響信号に対する 音源分離のための 非負値行列分解と半正定値テンソル分解

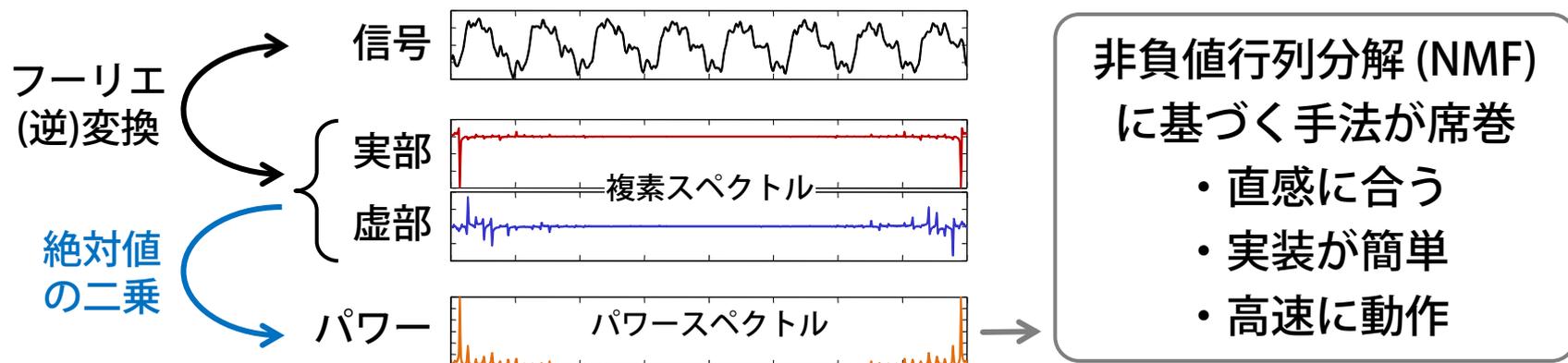
吉井 和佳 糸山 克寿 (京大)  
後藤 真孝 (産総研)

非ガウス性モノラル音響信号に対する  
音源分離のための  
非負値行列分解と半正定値テンソル分解

吉井 和佳 糸山 克寿 (京大)  
後藤 真孝 (産総研)

# 研究の背景

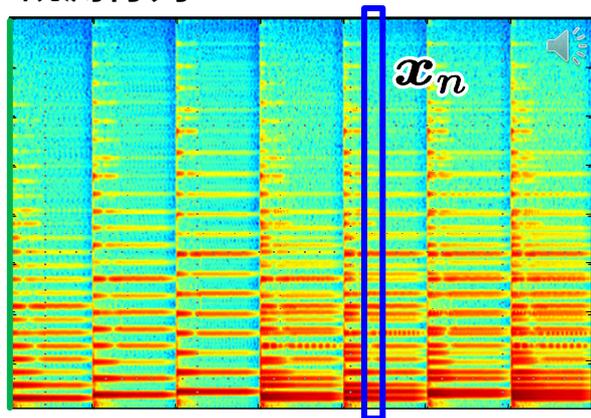
- モノラル音響信号の音源分離に関する研究がさかん
  - 応用例：音楽信号の高度な加工
    - 歌声と伴奏の分離 [Rafii2011, Huang2012, Yang2014, Ikemiya2015]
    - 楽器音イコライザ [吉井2006, 宮本2008, 糸山2008]
- スペクトル領域で音源分離を行うのが一般的
  - 位相を無視すると音の特徴 (調波構造・スパース性) がとらえやすくなる



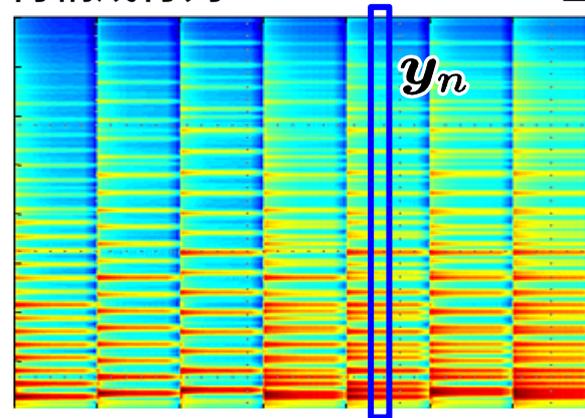
# 非負値行列分解 (NMF) [Lee 1999 / Févotte 2009]

- 「非負値ベクトル」を少数の「非負値ベクトル」の錐結合で表現

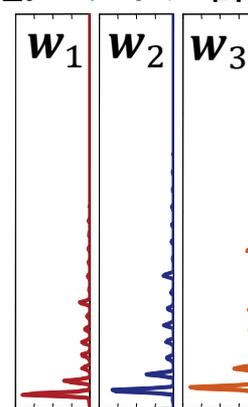
観測行列  $X$



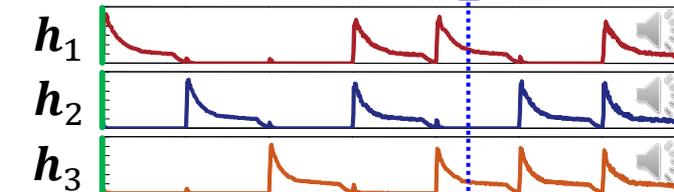
再構成行列  $Y = WH$



基底ベクトル群  $W$



$$\mathbf{x}_n \approx \mathbf{y}_n = \sum_{k=1}^K h_{kn} \mathbf{w}_k$$



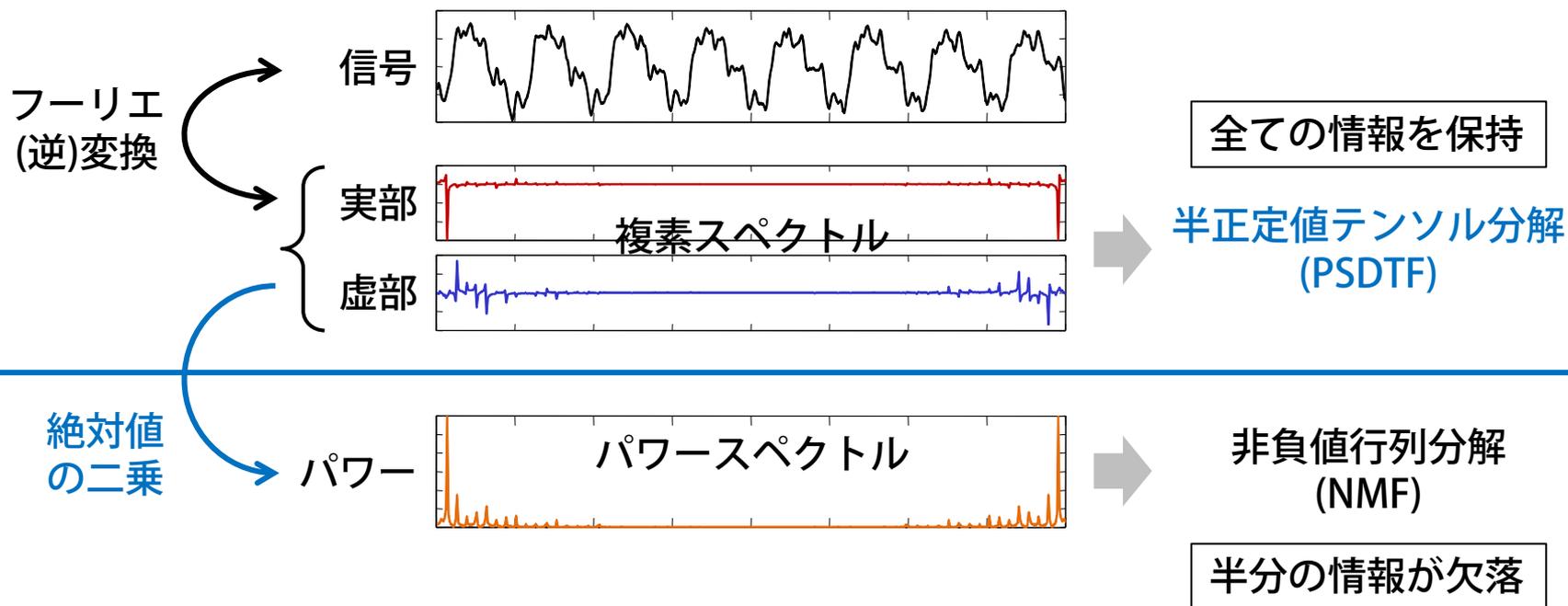
音量ベクトル群  $H$

音源分離に適したコスト関数：板倉・斎藤ダイバージェンス

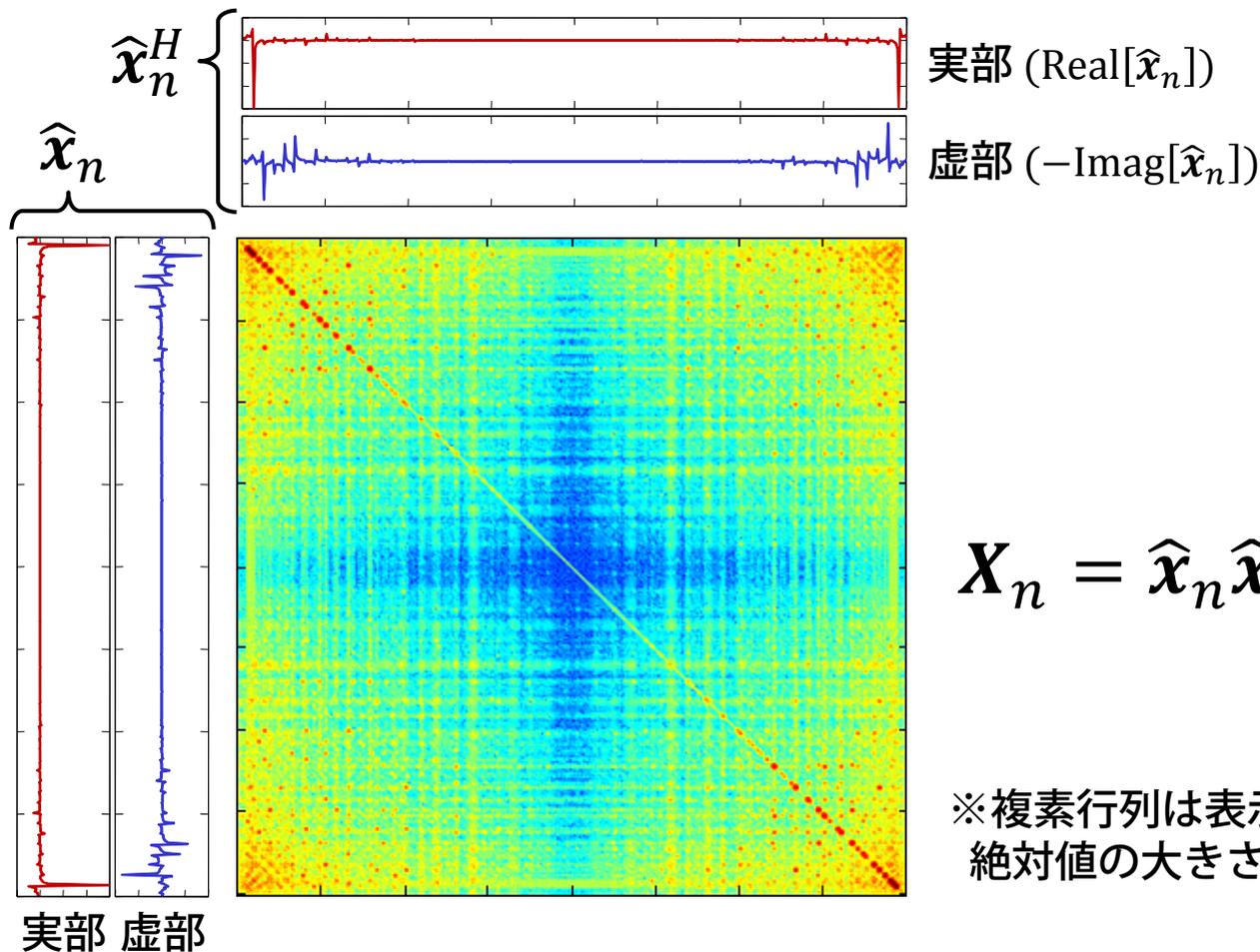
$$D_{IS}(\mathbf{x}_n | \mathbf{y}_n) = \sum_m (-\log x_{nm} y_{nm}^{-1} + x_{nm} y_{nm}^{-1} - 1)$$

# 研究の背景

- 複素スペクトログラムを直接分解できる方法論が必要
  - 入力音響信号のもつ情報 (パワー・位相) を全て適切に取り扱いたい
  - パワースペクトルにした時点で位相情報が失われる

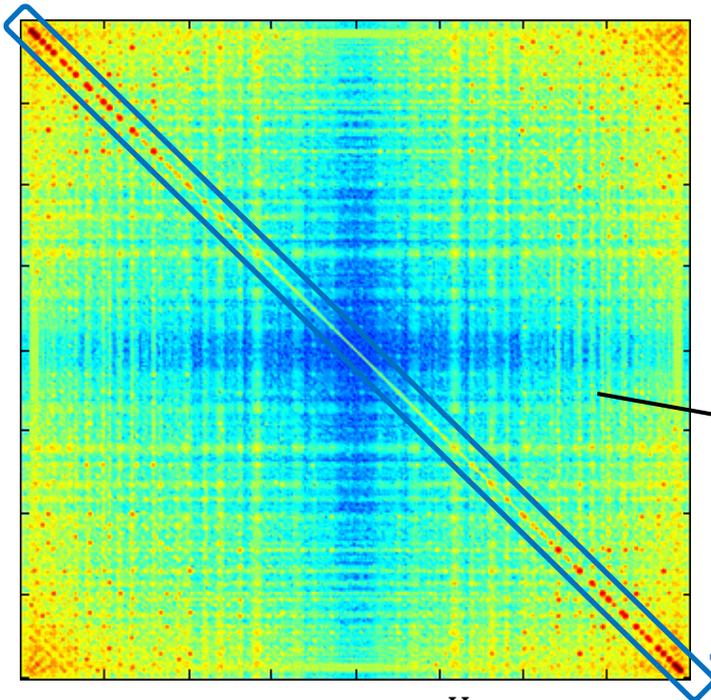


# 複素スペクトルの自己共分散



# 観測データの与え方

- 各フレームの複素スペクトルの自己共分散を計算 → 半正定値行列
  - 各フレームの複素スペクトルの絶対値の二乗を計算 → 非負値ベクトル



$$X_n = \hat{x}_n \hat{x}_n^H$$

ある行列 $X$ が半正定値行列であるとは

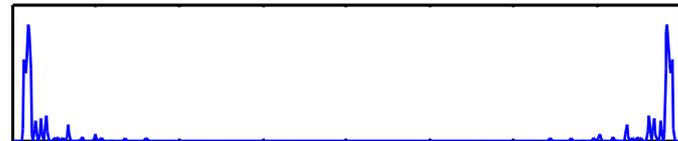
- すべての固有値が非負
- $X = AA^H$ となる行列 $A$ が存在

} 同値

半正定値性は非負値性の拡張概念

行列要素は任意のスカラ (負値や複素数もOK)

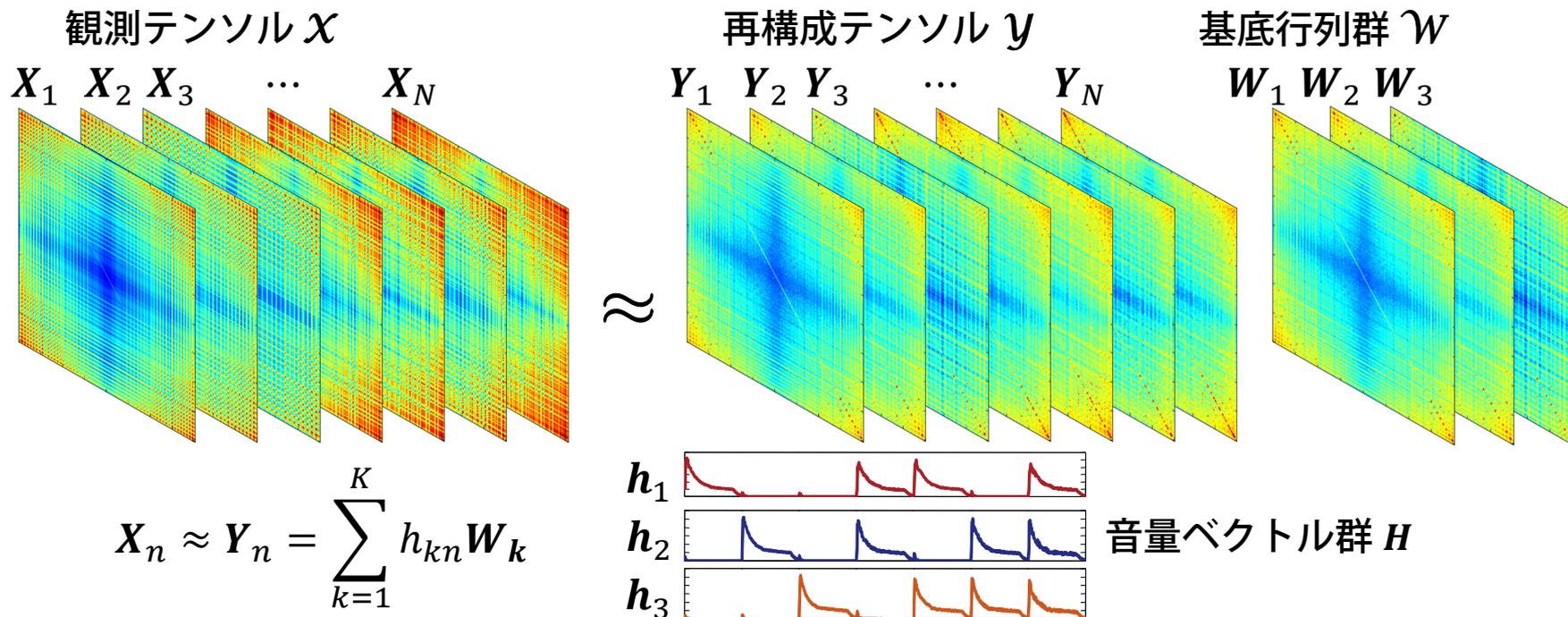
対角成分は常に非負 (パワースペクトル)



$$x_n = \hat{x}_n \odot \hat{x}_n$$

# 半正定値テンソル分解 (PSDTF) [Yoshii 2013]

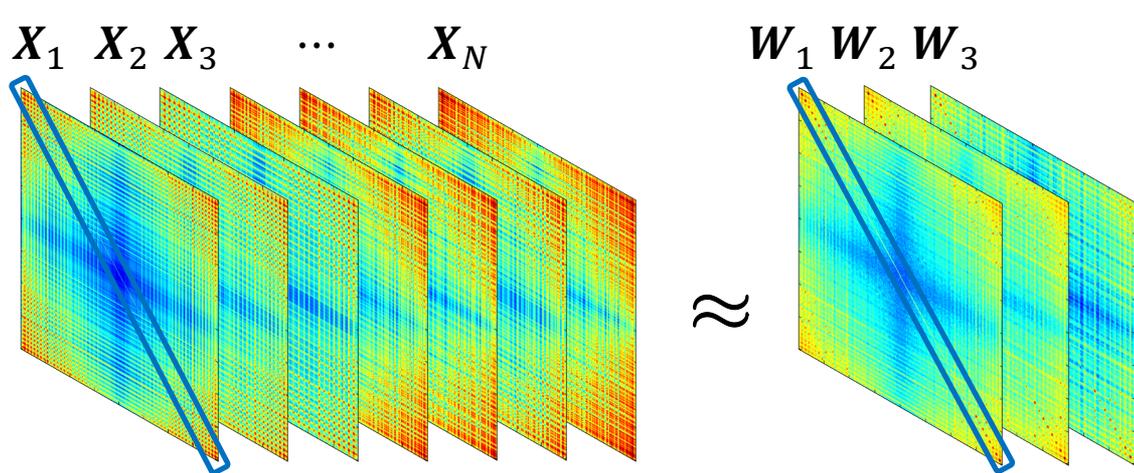
- 「半正定値行列」を少数の「半正定値行列」の錐結合で表現



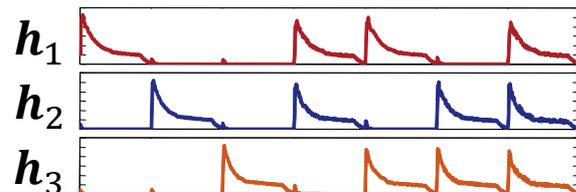
音源分離に適したコスト関数：Log-Determinantダイバージェンス

$$D_{LD}(X_n|Y_n) = -\log |X_n Y_n^{-1}| + X_n Y_n^{-1} - M$$

# PSDTF v.s. NMF



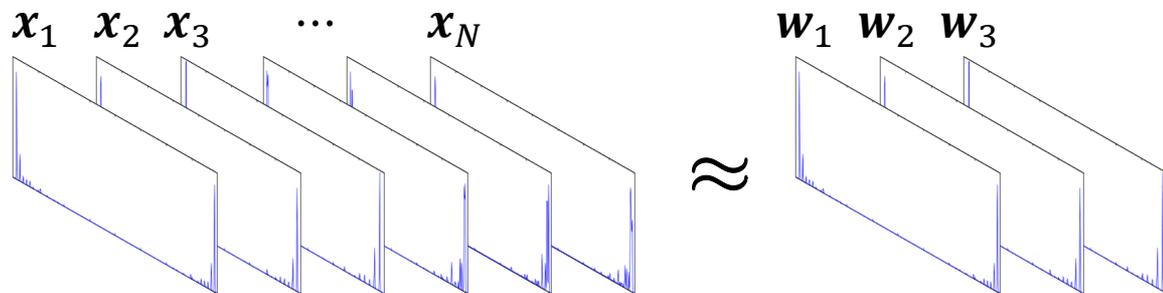
半正定値行列  $X_n \approx \sum_{k=1}^K h_{kn} W_k$  半正定値行列



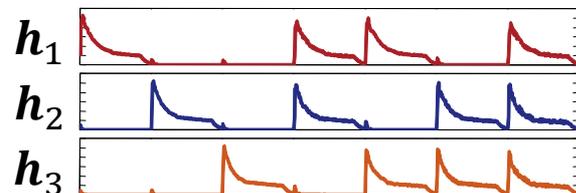
音源分離精度 : SDR 26.7 dB



対角成分に限定 (周波数ビン間の相関を無視)

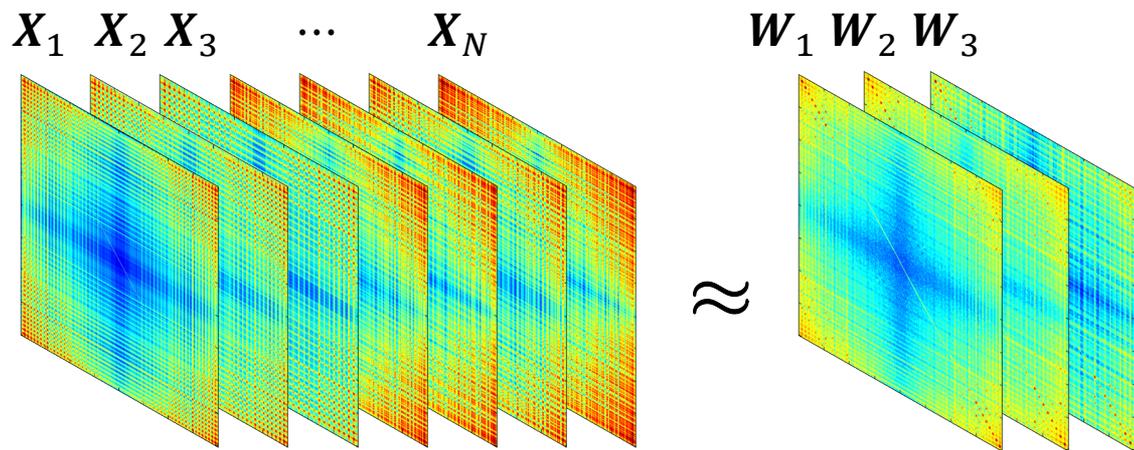


非負値ベクトル  $x_n \approx \sum_{k=1}^K h_{kn} w_k$  非負値ベクトル

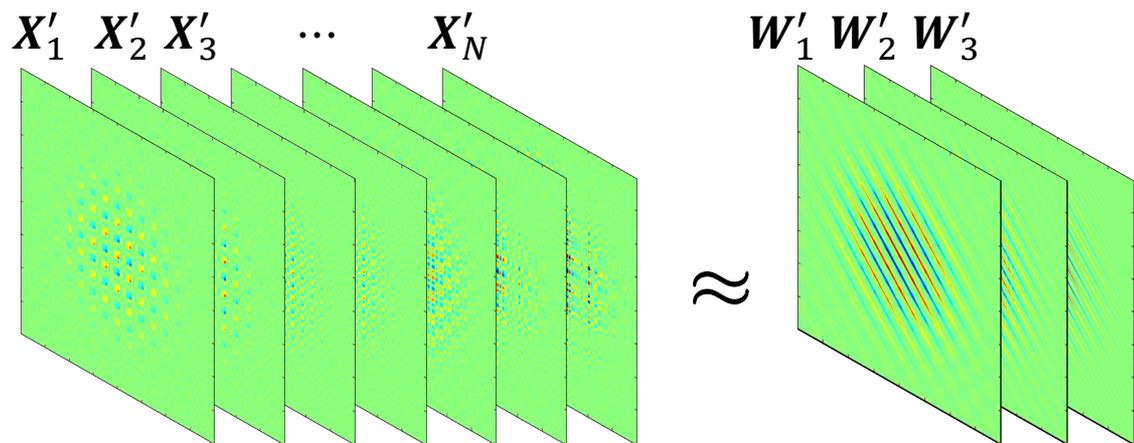


音源分離精度 : SDR 18.9 dB

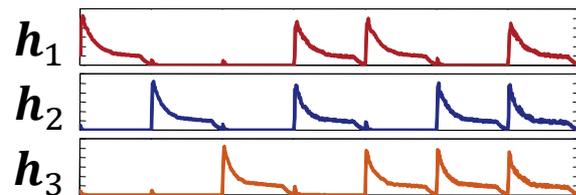
# 周波数領域PSDTF v.s. 時間領域PSDTF



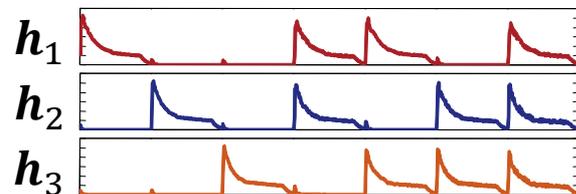
離散フーリエ変換行列 $F$ で線形変換



$$\text{半正定値行列 } (FX'_n F^H) \approx \sum_{k=1}^K \text{半正定値行列 } h_{kn} (FW'_k F^H)$$



$$\text{半正定値行列 } X'_n \approx \sum_{k=1}^K \text{半正定値行列 } h_{kn} W'_k$$



# 非ガウス性モノラル音響信号に対する 音源分離のための 非負値行列分解と半正定値テンソル分解

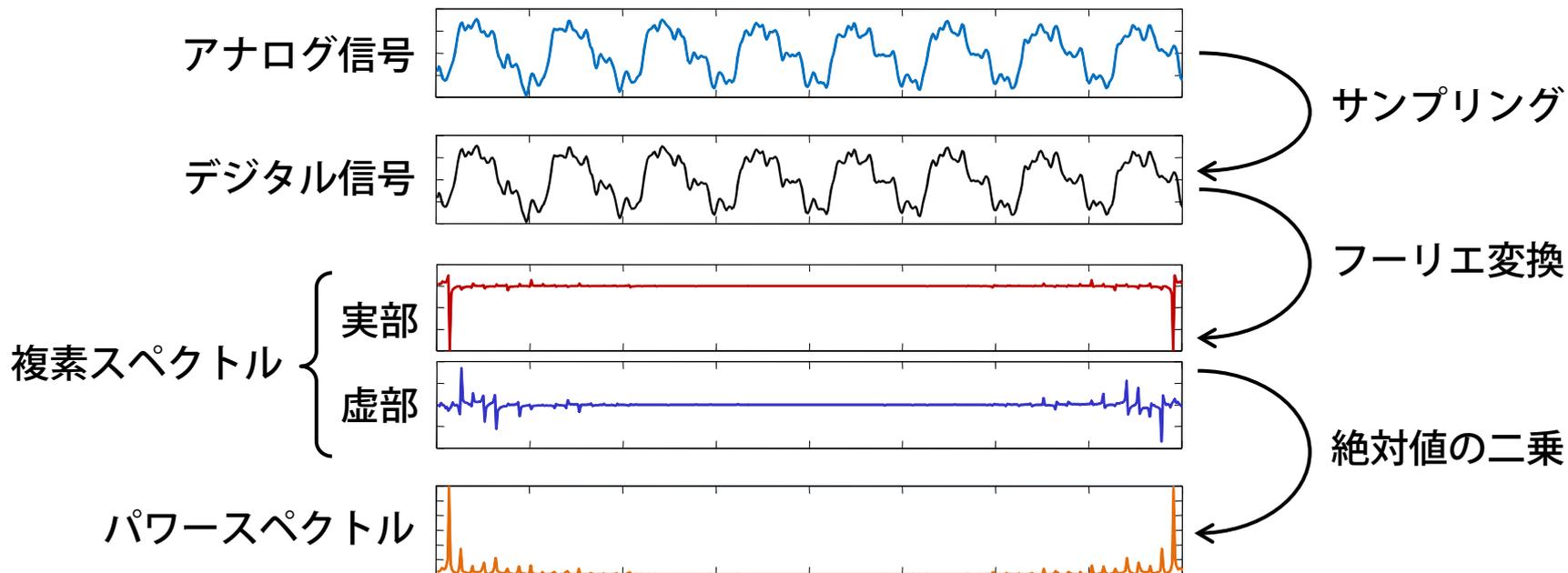
吉井 和佳 糸山 克寿 (京大)  
後藤 真孝 (産総研)

# 非ガウス性モノラル音響信号に対する 音源分離のための 非負値行列分解と半正定値テンソル分解

吉井 和佳 糸山 克寿 (京大)  
後藤 真孝 (産総研)

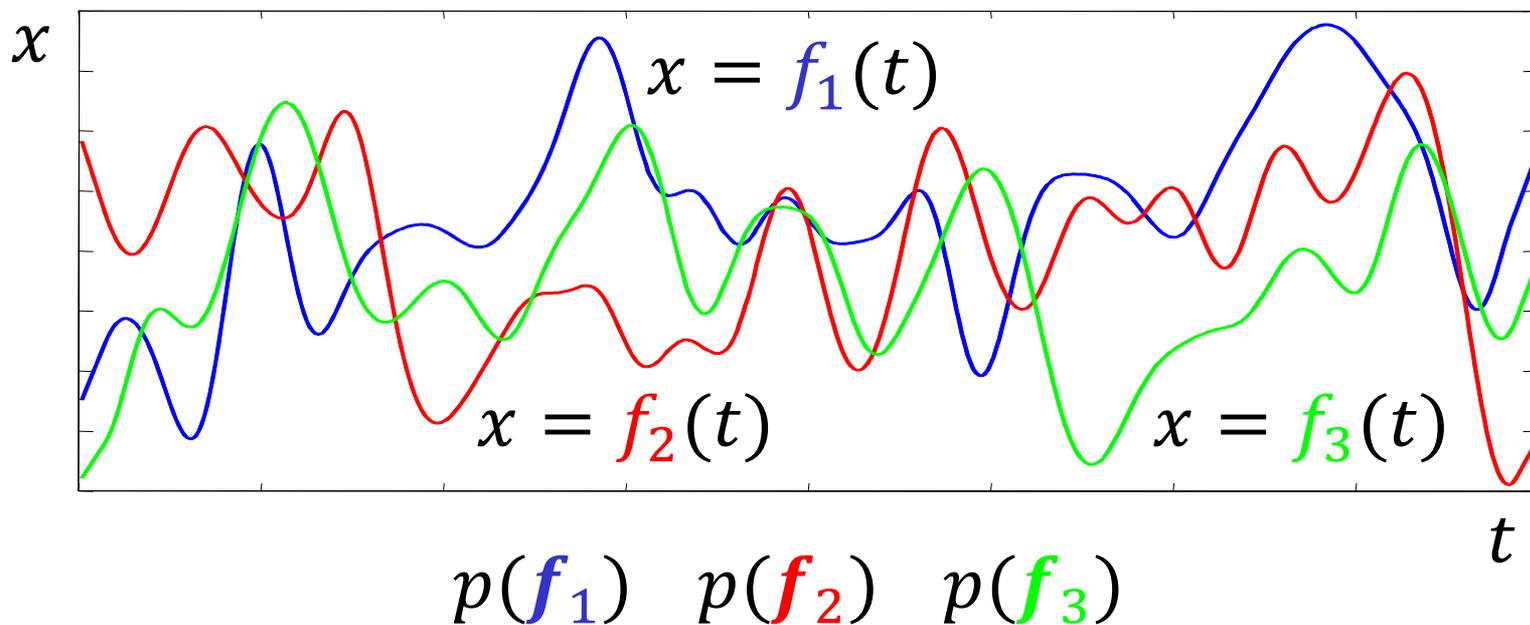
# 研究の動機

- 音響信号は連続信号 → 連続信号に対する確率モデルが必要
  - 空気振動は連続時間上の物理現象
  - 「サンプリング」によって離散信号が得られる



# ガウス過程 (GP)

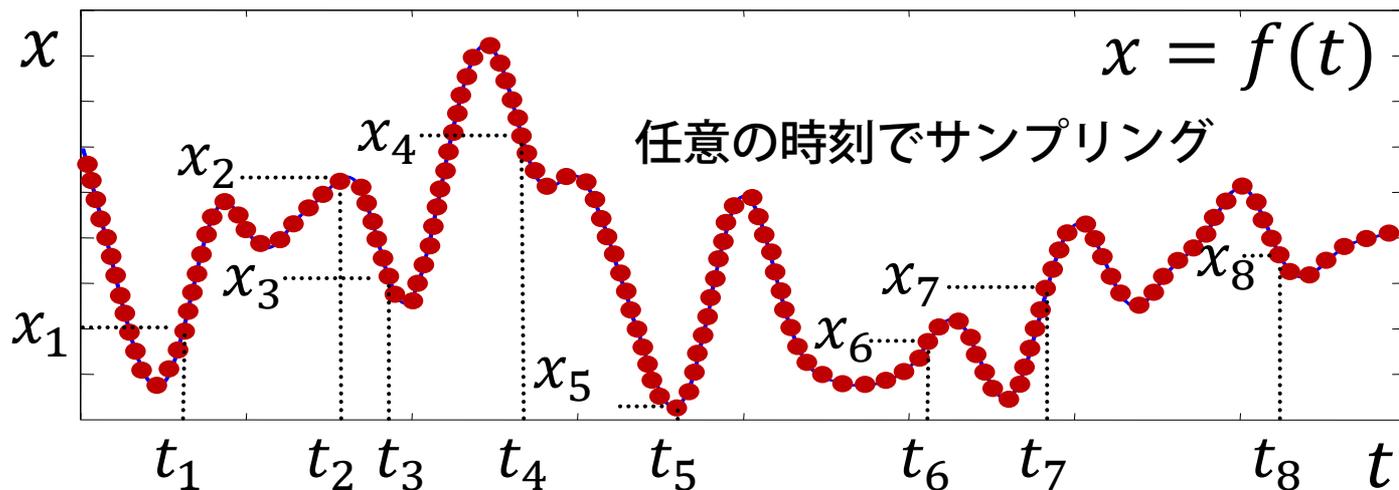
- 連続時間上の確率過程 (無限次元の確率分布) のひとつ
  - 連続関数  $f$  (無限次元のベクトル) に対する確率分布  $p(f)$
  - 任意の音響信号のもっともらしさを計算可能



# ガウス過程とガウス分布

- ガウス過程  $\Leftrightarrow$  無限次元のガウス分布

– 連続関数  $f \Leftrightarrow$  無限個の出力値  $x = [x_1, x_2, \dots, x_\infty]^T$



関数  $f$   $\begin{cases} \rightarrow t = [t_1, t_2, \dots, t_\infty]^T \\ \rightarrow x = [x_1, x_2, \dots, x_\infty]^T \end{cases}$



無限次元のガウス分布

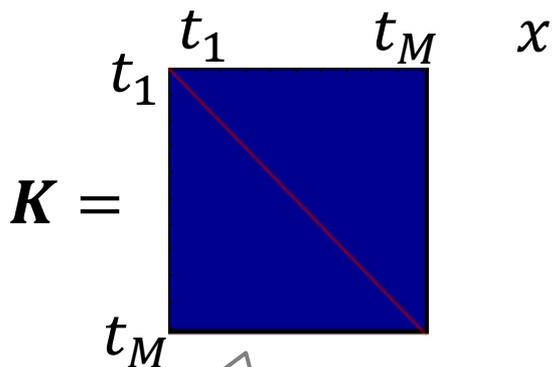
$$x \sim \text{GP}(0, K)$$

任意の周辺分布はガウス分布

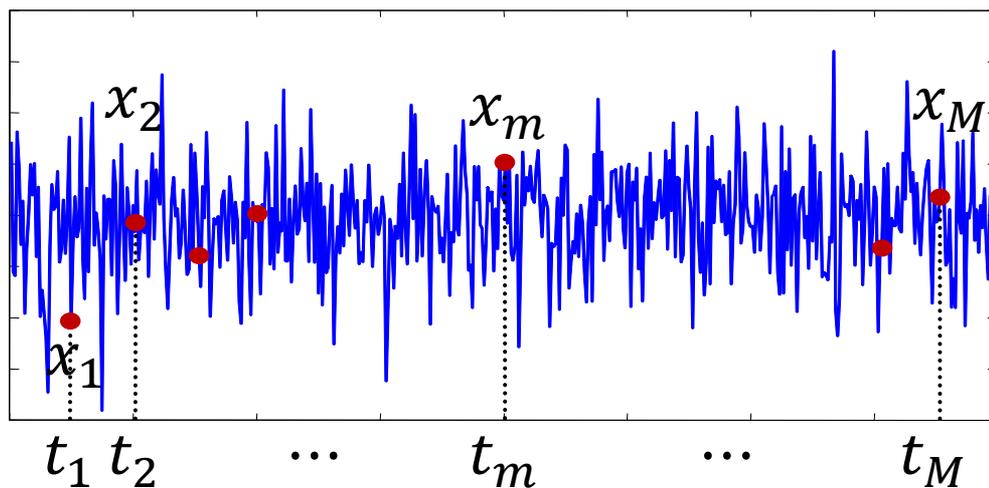
# Deltaカーネル

- 定常な白色雑音を生成
  - 周辺分布  $\mathbf{x} \sim N(\mathbf{0}, \mathbf{K})$

関数  $f$   $\begin{cases} \mathbf{t} = [t_1, t_2, \dots, t_M]^T \\ \mathbf{x} = [x_1, x_2, \dots, x_M]^T \end{cases}$



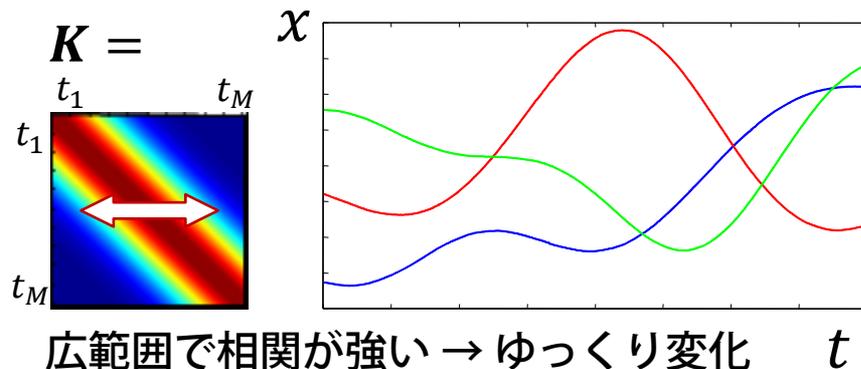
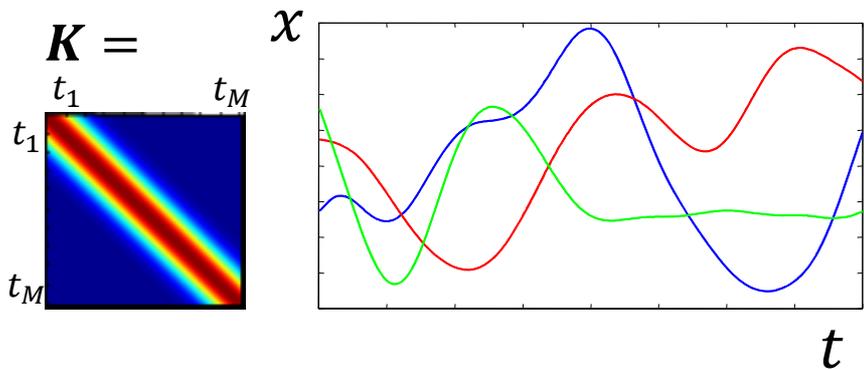
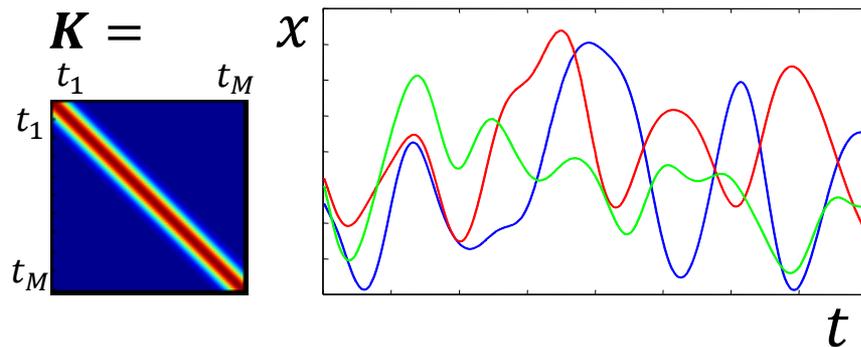
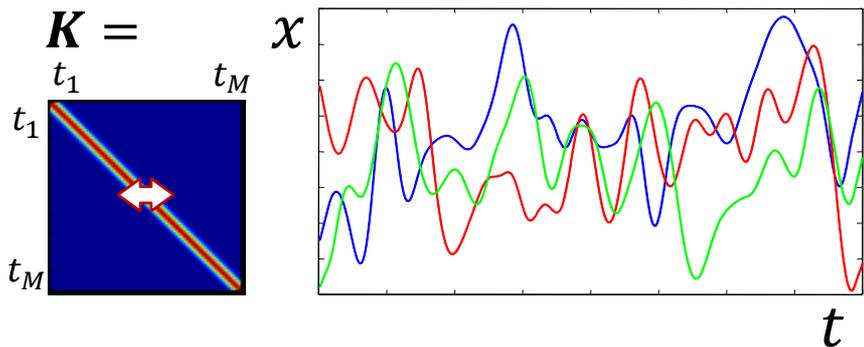
各  $x_m$  は独立な  
ガウス分布に従う



# Squared Exponentialカーネル

- なめらかな連続関数を生成
  - 周辺分布  $\mathbf{x} \sim N(\mathbf{0}, \mathbf{K})$

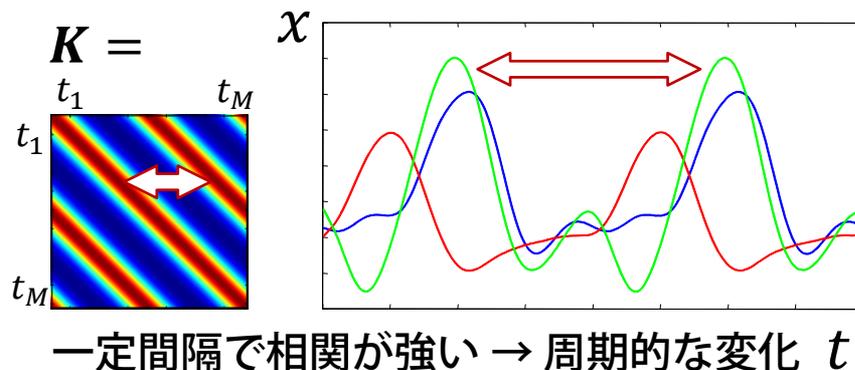
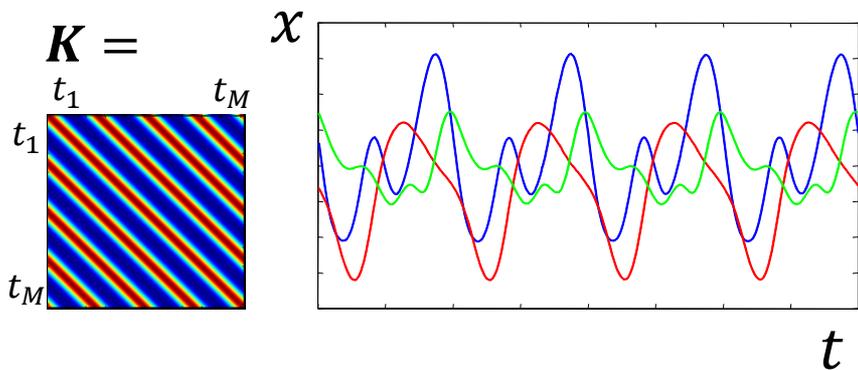
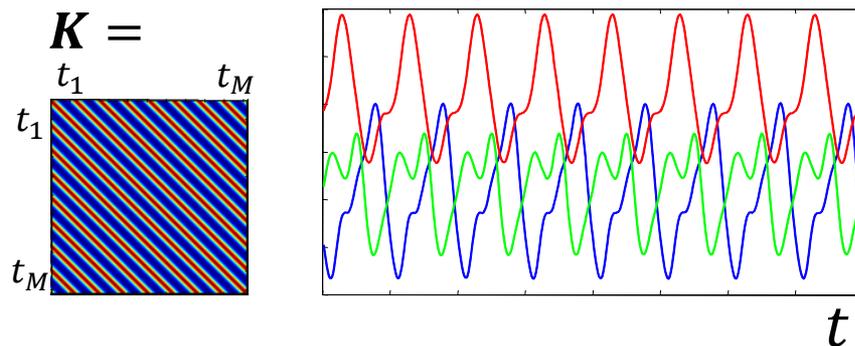
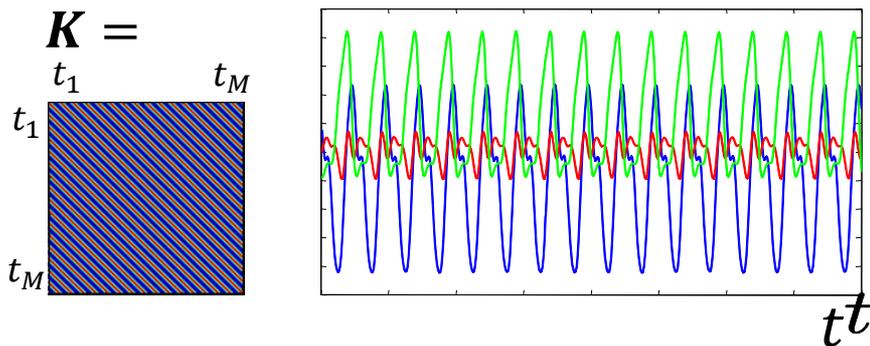
関数  $f$   $\begin{cases} \mathbf{t} = [t_1, t_2, \dots, t_M]^T \\ \mathbf{x} = [x_1, x_2, \dots, x_M]^T \end{cases}$



# 周期カーネル

- 周期的な連続関数を生成
  - 周辺分布  $\mathbf{x} \sim N(\mathbf{0}, \mathbf{K})$

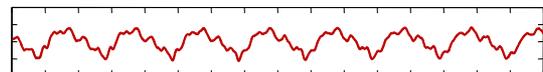
関数  $f$   $\begin{cases} \rightarrow \mathbf{t} = [t_1, t_2, \dots, t_M]^T \\ \rightarrow \mathbf{x} = [x_1, x_2, \dots, x_M]^T \end{cases}$



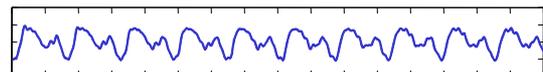
# 混合音の生成と分離

- 生成：ガウス変数の和 → ガウス変数

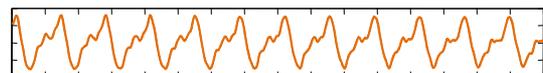
再生性をもつ



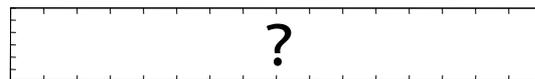
$$\mathbf{x}_1 \sim N(\mathbf{0}, \mathbf{K}_1)$$



$$\mathbf{x}_2 \sim N(\mathbf{0}, \mathbf{K}_2)$$



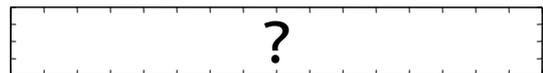
$$\mathbf{x}_3 \sim N(\mathbf{0}, \mathbf{K}_3)$$



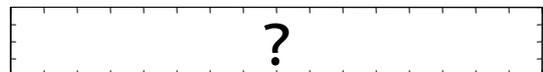
$$\underbrace{\mathbf{x}_1 + \mathbf{x}_2 + \mathbf{x}_3}_{\mathbf{x}} \sim N(\mathbf{0}, \underbrace{\mathbf{K}_1 + \mathbf{K}_2 + \mathbf{K}_3}_{\mathbf{K}})$$

- 分離：ガウス変数の分解 → ガウス変数

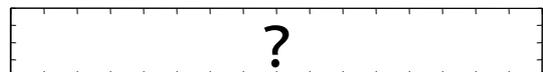
条件付き期待値が求まる



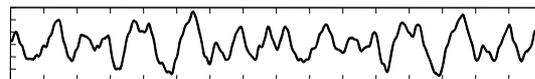
$$\mathbf{x}_1 \sim N(\mathbf{0}, \mathbf{K}_1)$$



$$\mathbf{x}_2 \sim N(\mathbf{0}, \mathbf{K}_2)$$



$$\mathbf{x}_3 \sim N(\mathbf{0}, \mathbf{K}_3)$$



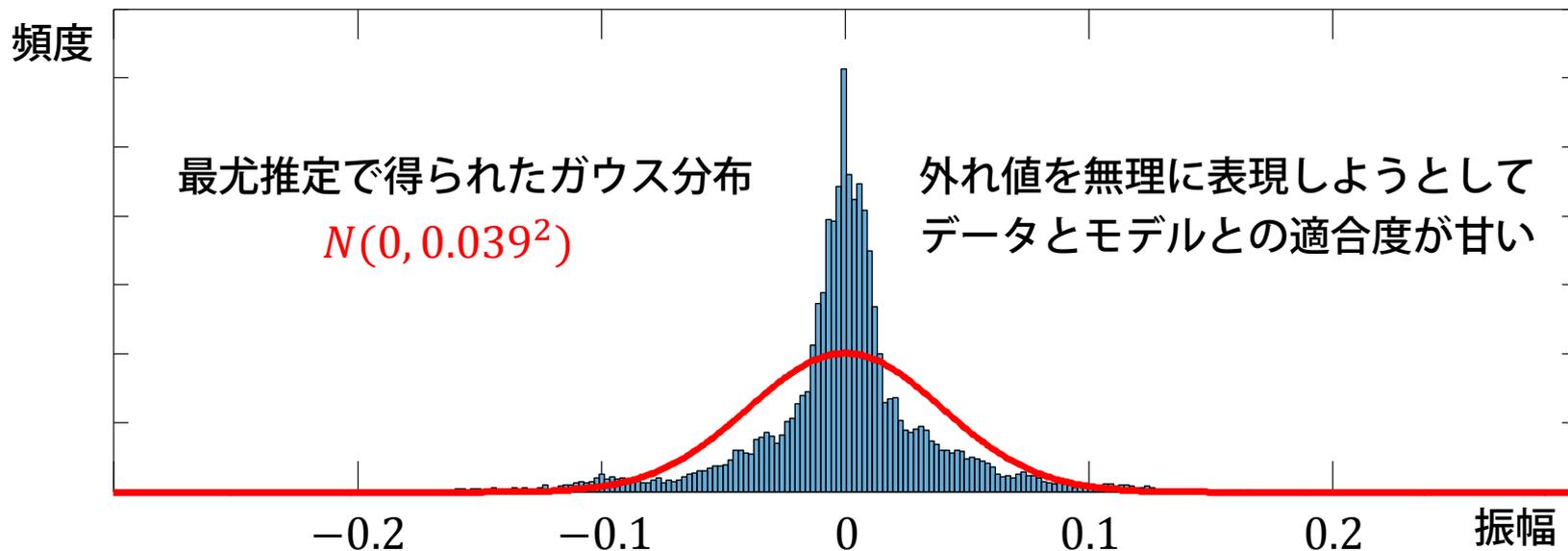
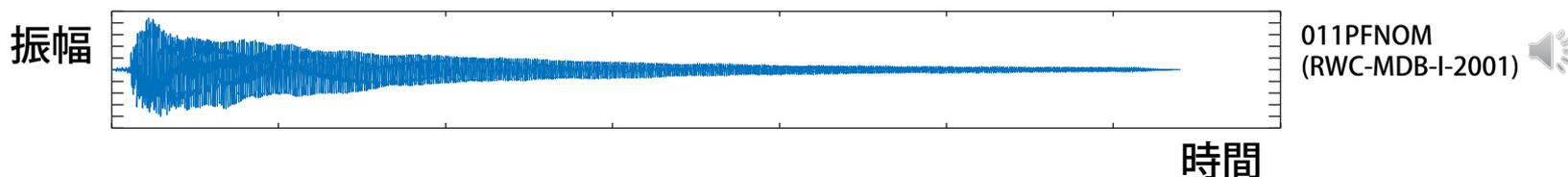
$$\mathbf{x}_1 + \mathbf{x}_2 + \mathbf{x}_3 \sim N(\mathbf{0}, \mathbf{K}_1 + \mathbf{K}_2 + \mathbf{K}_3)$$

ウィナーフィルタによる音源分離

$$E[\mathbf{x}_1|\mathbf{x}] \sim \mathbf{K}_1\mathbf{K}^{-1}\mathbf{x} \quad E[\mathbf{x}_2|\mathbf{x}] \sim \mathbf{K}_2\mathbf{K}^{-1}\mathbf{x} \quad E[\mathbf{x}_3|\mathbf{x}] \sim \mathbf{K}_3\mathbf{K}^{-1}\mathbf{x}$$

# ガウス分布の限界

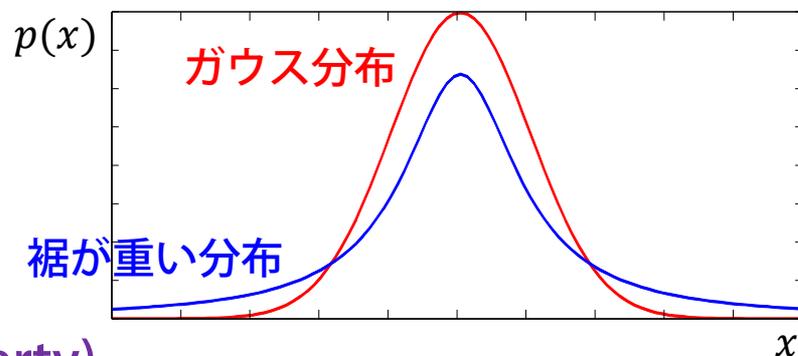
- 実際の音源信号はガウス分布には従わないことが多い
  - ガウス分布では値の散らばり (はずれ値含む) を十分に表現できない



# どのような確率分布を用いるべきか

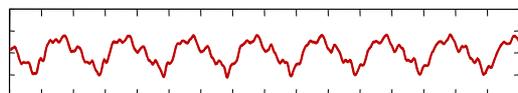
- 裾が重い (heavy-tailed)

- 裾が指数関数的には減衰しない
- 外れ値に頑健なモデル化が可能

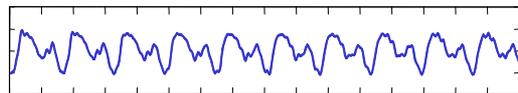


- 再生性をもつ (reproductive property)

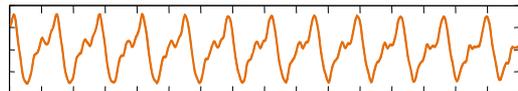
- 確率変数の和が同じクラスの確率分布に従う
- 音源信号の加法性を成立させるためには不可欠



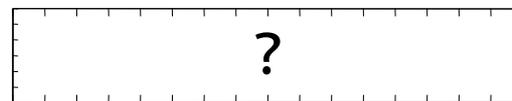
$$x_1 \sim N(\mathbf{0}, K_1)$$



$$x_2 \sim N(\mathbf{0}, K_2)$$



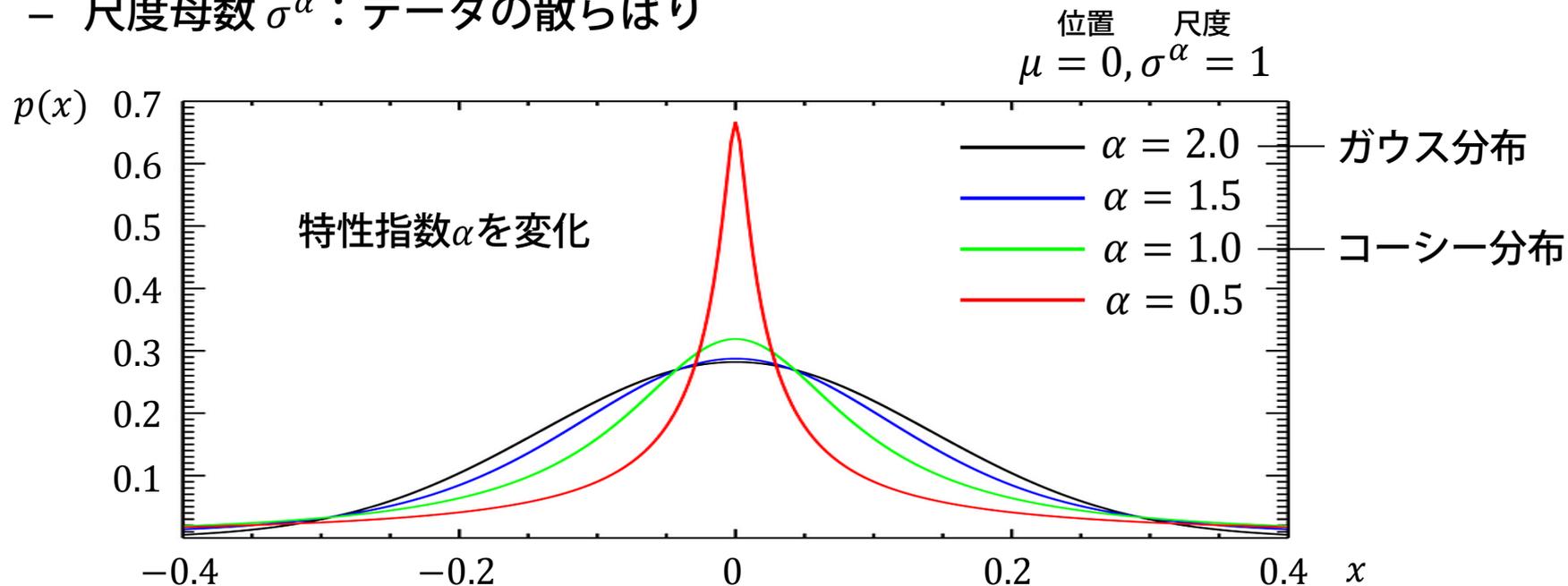
$$x_3 \sim N(\mathbf{0}, K_3)$$



$$x_1 + x_2 + x_3 \sim N(\mathbf{0}, K_1 + K_2 + K_3)$$

# 対称 $\alpha$ 安定分布

- 3個のパラメータを持つ左右対称な確率分布  $S_\alpha(x|\mu, \sigma^\alpha)$ 
  - 特性指数  $\alpha$  :  $S_\alpha$ 分布を特徴づける最も重要な量 ( $0 < \alpha \leq 2$ )
  - 位置母数  $\mu$  : 本研究では0 (音響信号は0を中心に振動)
  - 尺度母数  $\sigma^\alpha$  : データの散らばり



# S $\alpha$ S分布の性質

- 理論的にはS $\alpha$ S分布が音源分離に最も適している [Liutkus 2015]

- 裾が重い  $\rightarrow$  音源信号に含まれる外れ値に対応
- 再生性をもつ  $\rightarrow$  音源信号の加法性が成立

非ガウス信号分離のための  
一般化ウィナーフィルタを提案

「安定」分布と呼ばれる由来

$$X_1, X_2 \sim S_\alpha(0, \sigma^\alpha) \quad \Rightarrow \quad AX_1 + BX_2 \sim S_\alpha(0, (A^\alpha + B^\alpha)\sigma^\alpha)$$

$(A, B > 0)$  やはり安定分布!

確率変数の加法性

$$\begin{cases} X_1 \sim S_\alpha(0, \sigma_1^\alpha) \\ X_2 \sim S_\alpha(0, \sigma_2^\alpha) \end{cases} \quad \Rightarrow \quad AX_1 + BX_2 \sim S_\alpha(0, A^\alpha \sigma_1^\alpha + B^\alpha \sigma_2^\alpha)$$

中心極限定理の例外!

$\alpha = 1$  (コーシー分布) なら振幅領域で  
 $\alpha = 2$  (ガウス分布) ならパワー領域で **加法性が成立**

# S $\alpha$ S分布の性質

- 一般の $0 < \alpha \leq 2$ に対する確率密度関数 $p(x)$ は解析的に計算不可能
  - ガウス分布 ( $\alpha = 2$ ) とコーシー分布 ( $\alpha = 1$ ) のみ解析的に表現可能
  - ただし、特性関数 $\psi(t)$ は簡潔に表現可能

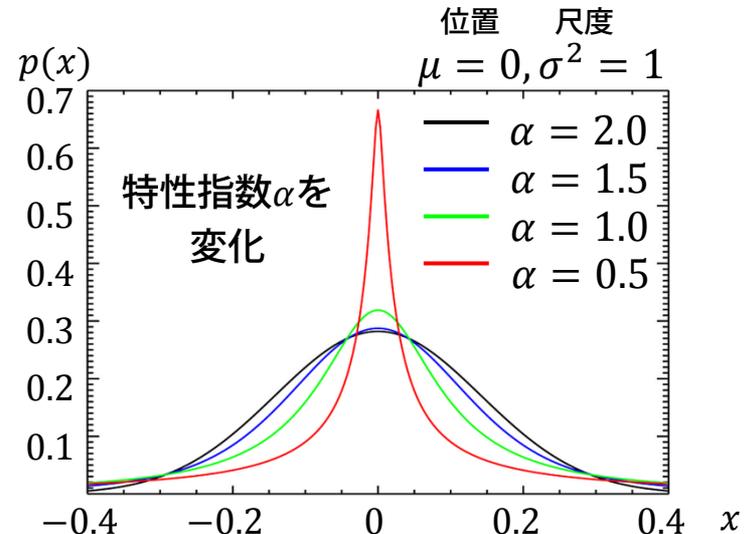
$$\psi(t) = E_x[e^{itx}] = \int_{-\infty}^{\infty} e^{itx} p(x) dx = \exp(i\mu t - \sigma^\alpha |t|^\alpha)$$

$p(x)$ のフーリエ変換



次善の策として用いるべき確率分布の要件

- 裾が重い
- ガウス/コーシー分布を特殊形に含む
- 確率密度関数が解析的に表現可能



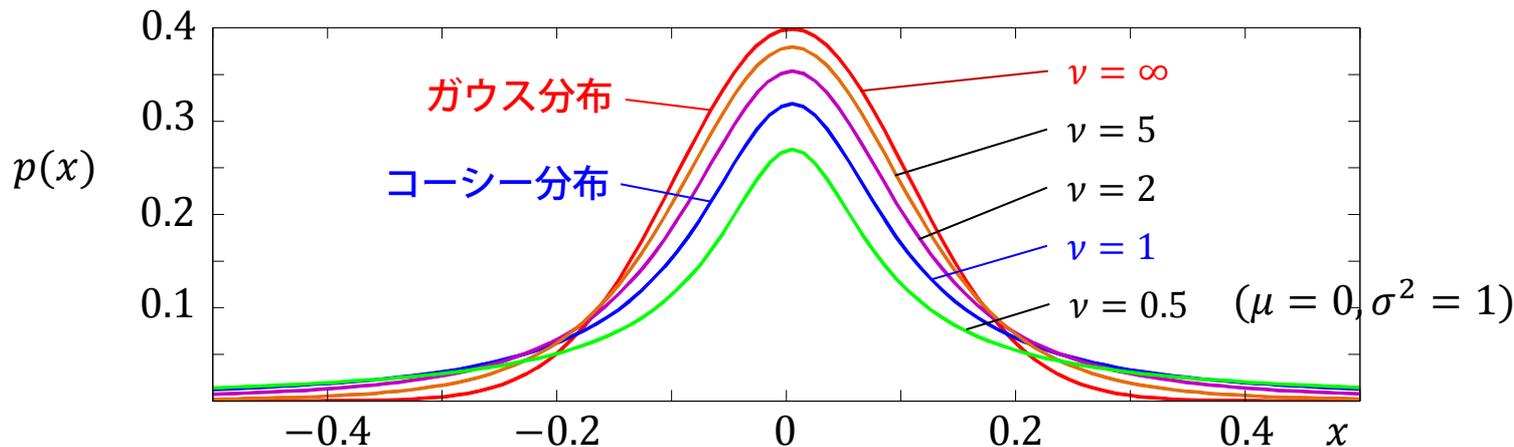
# t分布

- ガウススケール混合分布の一種 (Gaussian scale mixture: GSM)
  - 裾が重い・コーシー分布・ガウス分布を特殊形に含む

$$p(x|\mu, \sigma^2, \nu) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\Gamma\left(\frac{\nu}{2}\right)} \frac{1}{\sqrt{\pi\nu\sigma^2}} \left(1 + \frac{(x-\mu)^2}{\sigma^2}\right)^{-\frac{\nu+1}{2}} \rightarrow N(x|\mu, \sigma^2) = \frac{1}{\sqrt{\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{\sigma^2}\right)$$

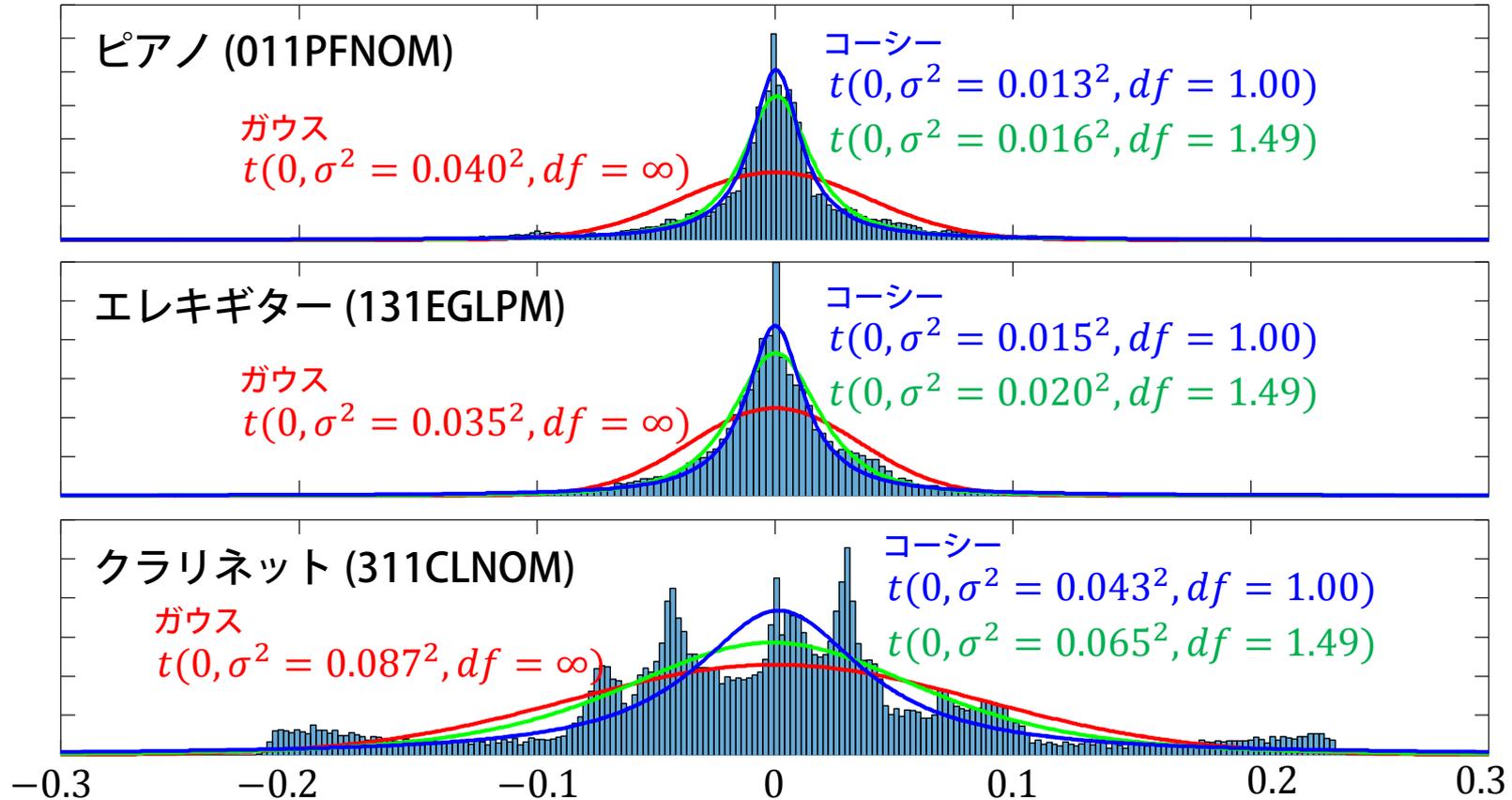
自由度

自由度  $\nu \rightarrow \infty$  でガウス分布に一致



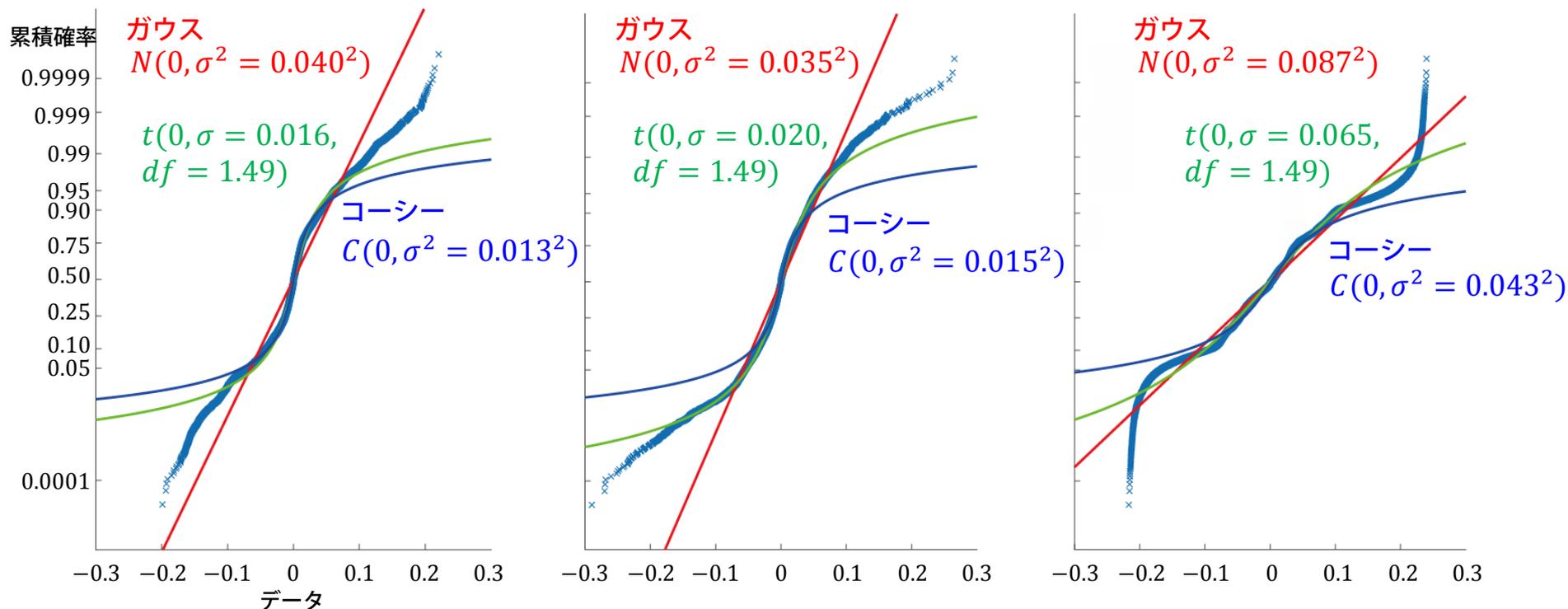
# t分布による音響信号のモデル化

- t分布はガウス分布より当てはまりがよい



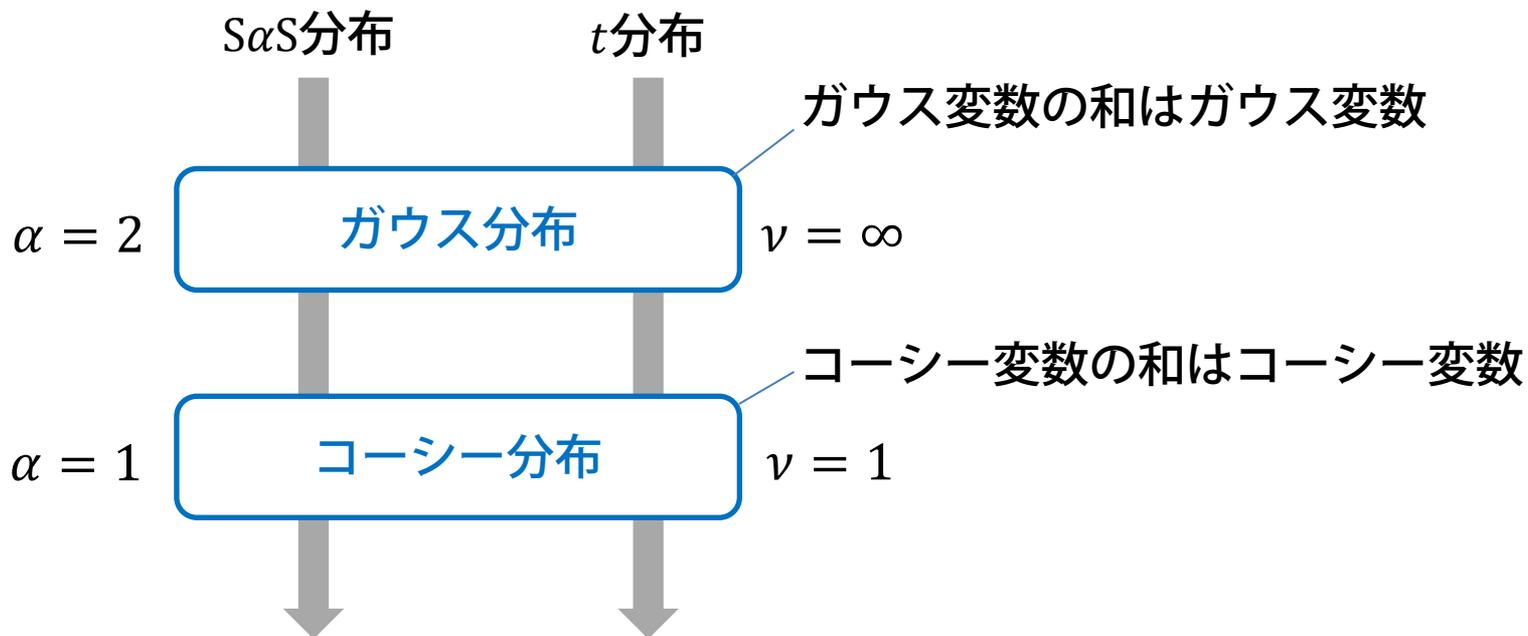
# t分布による音響信号のモデル化

- t分布はガウス分布より当てはまりがよい
  - 中心部の当てはまりは良いが、裾の厚みを大きく見積もりすぎる傾向



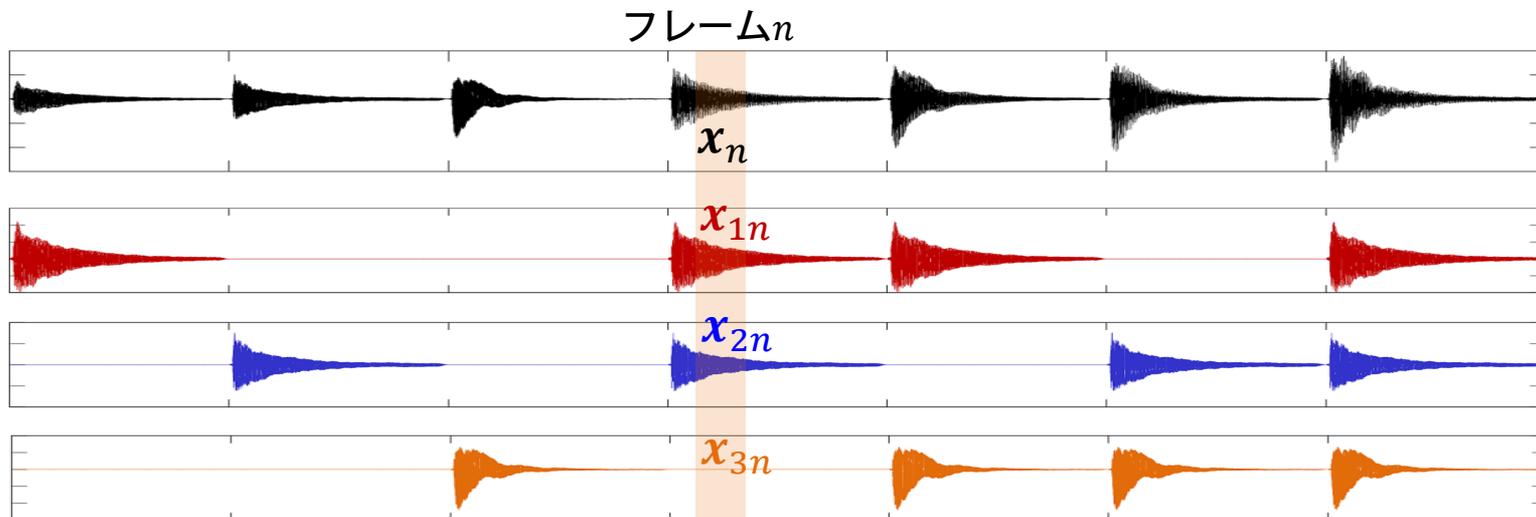
# S $\alpha$ S分布とt分布

- いずれもガウス分布とコーシー分布を特殊形として含む
  - 本来はS $\alpha$ S分布が適切 [Liutkus 2015]
  - 確率密度関数が解析的に計算可能なt分布に着目



# 混合音の生成モデル

- ガウス変数 (局所信号) の足し合わせはガウス変数



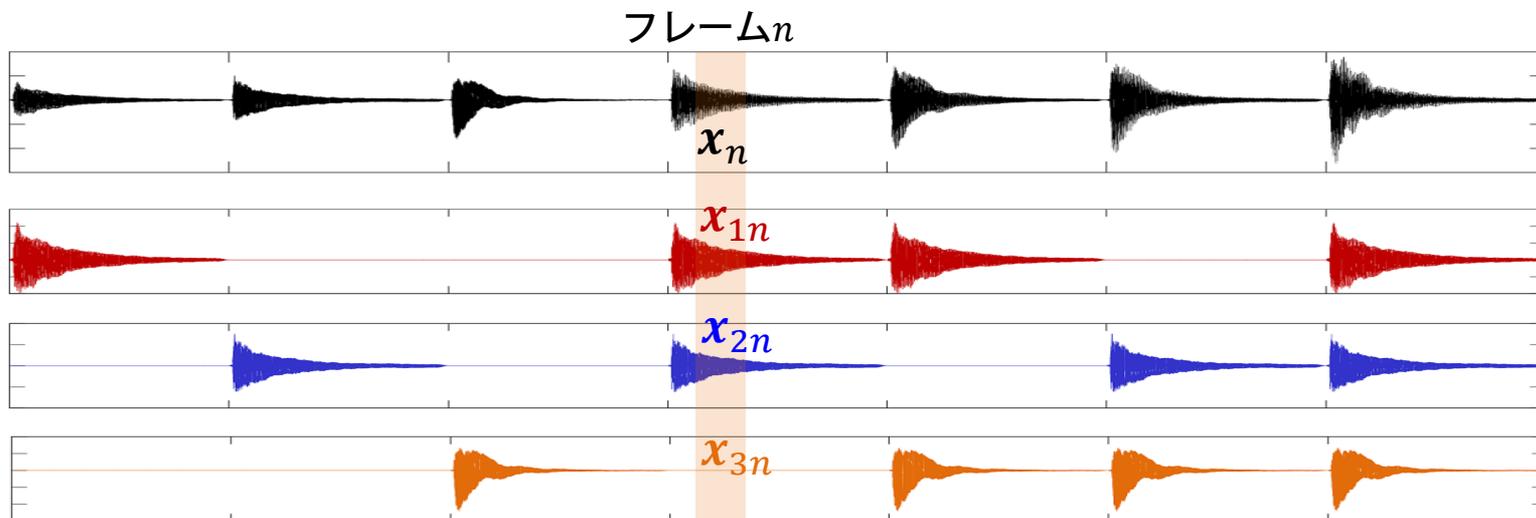
$$\left. \begin{aligned} x_{1n} &\sim N(\mathbf{0}, h_{1n} \mathbf{W}_1) \\ x_{2n} &\sim N(\mathbf{0}, h_{2n} \mathbf{W}_2) \\ x_{3n} &\sim N(\mathbf{0}, h_{3n} \mathbf{W}_3) \end{aligned} \right\} x_n \sim N\left(\mathbf{0}, \sum_{k=1}^K h_{kn} \mathbf{W}_k\right)$$

$x_n x_n^T \geq 0$  が与えられた時に  
 $h_{kn} \geq 0$  と  $\mathbf{W}_k \geq 0$  を求めたい

ガウス過程に基づく PSDTF (時間領域)

# 混合音の生成モデル

- コーシー変数 (局所信号) の足し合わせはコーシー変数



$$\left. \begin{aligned} \mathbf{x}_{1n} &\sim C(\mathbf{0}, h_{1n} \mathbf{W}_1) \\ \mathbf{x}_{2n} &\sim C(\mathbf{0}, h_{2n} \mathbf{W}_2) \\ \mathbf{x}_{3n} &\sim C(\mathbf{0}, h_{3n} \mathbf{W}_3) \end{aligned} \right\} \mathbf{x}_n \sim C\left(\mathbf{0}, \sum_{k=1}^K h_{kn} \mathbf{W}_k\right)$$

$\mathbf{x}_n \mathbf{x}_n^T \geq 0$  が与えられた時に  
 $h_{kn} \geq 0$  と  $\mathbf{W}_k \geq 0$  を求めたい

コーシー過程に基づく PSDTF (時間領域)

# $t$ -PSDTFの定式化

- 尤度関数としてガウス過程の代わりに $t$ 過程を利用

- 加法性が成立する場合は限られる

- $\nu = 1$  : コーシー過程に基づくPSDTF
- $\nu = \infty$  : ガウス過程に基づくPSDTF

$$\mathbf{x}_n \sim t \left( \mathbf{0}, \sum_{k=1}^K h_{kn} \mathbf{W}_k, \nu \right)$$

- 他の $\nu$ でも実用上うまくいく可能性はある

- パラメータ $H, \mathcal{W}$ の更新則は補助関数を用いて導出

- ガウス過程に基づくPSDTFの更新式を含む

- 周波数領域での等価なモデル化も可能

離散フーリエ  
変換行列 $F$ による  
線形変換

$$\mathbf{F}\mathbf{x}_n = \hat{\mathbf{x}}_n \sim t_c \left( \mathbf{0}, \sum_{k=1}^K h_{kn} \mathbf{F}\mathbf{W}_k \mathbf{F}^H, \nu \right)$$

複素スペクトル

# t-NMFの定式化

- 複素 $t$ 分布に基づく $t$ -NMF (周波数領域)

IS-NMF  $\hat{x}_{nm} \sim N_c \left( 0, \sum_{k=1}^K h_{kn} w_{km}, \nu \right)$   $\rightarrow$   $t$ -NMF  $\hat{x}_{nm} \sim t_c \left( 0, \sum_{k=1}^K h_{kn} w_{km}, \nu \right)$

複素スペクトルの値  
(フレーム $n$ ・ビン $m$ )

更新式  $w_{km} \leftarrow w_{km} \sqrt{\frac{\sum_n \frac{(\pi_n x_{nm}) h_{kn}}{y_{nm}^2}}{\sum_n \frac{h_{kn}}{y_{nm}}}}$   $h_{kn} \leftarrow h_{kn} \sqrt{\frac{\sum_m \frac{(\pi_n x_{nm}) w_{km}}{y_{nm}^2}}{\sum_m \frac{w_{km}}{y_{nm}}}}$

$\pi_n \rightarrow 1$  ( $\nu \rightarrow \infty$ )でIS-NMFの更新式に帰着

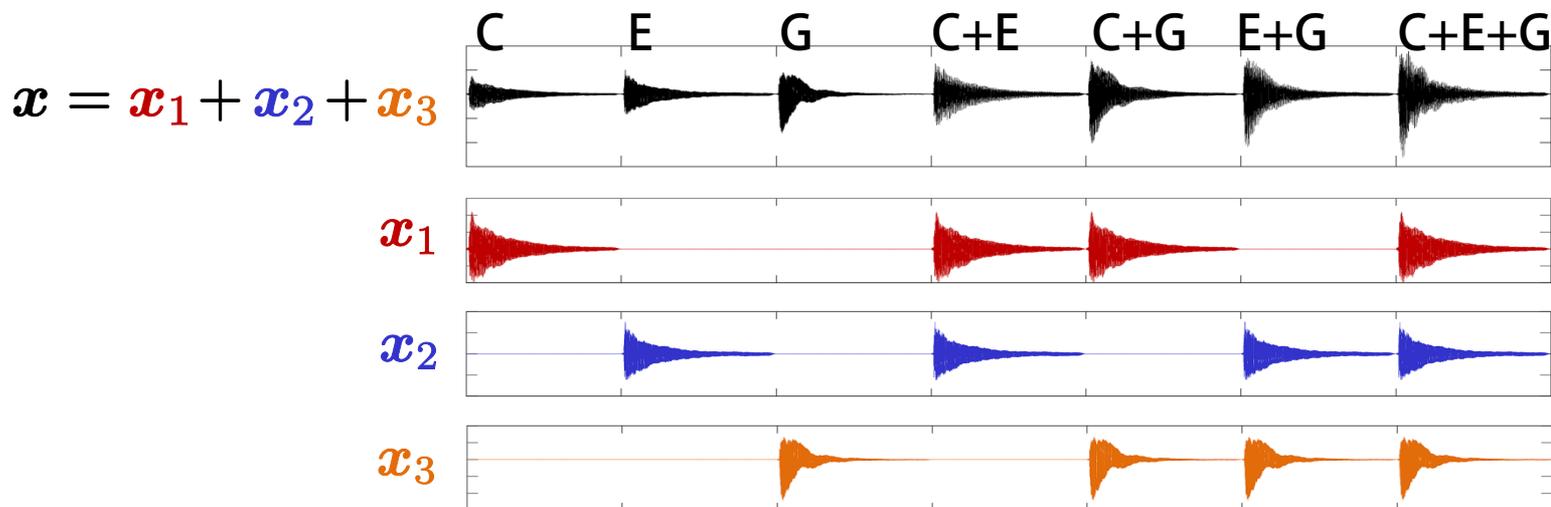
$$\pi_n x_{nm} = \left( \frac{2}{2 + \nu} y_{nm}^{-1} + \frac{\nu}{2 + \nu} x_{nm}^{-1} \right)^{-1}$$

$\rightarrow x_{nm}$  ( $\nu \rightarrow \infty$ )

観測データ $x_{nm}$ と再構成データ $y_{nm}$ を  
 $\nu: 2$ で調和平均をとったものを観測データと  
みなしてIS-NMF (過学習の抑制効果?)

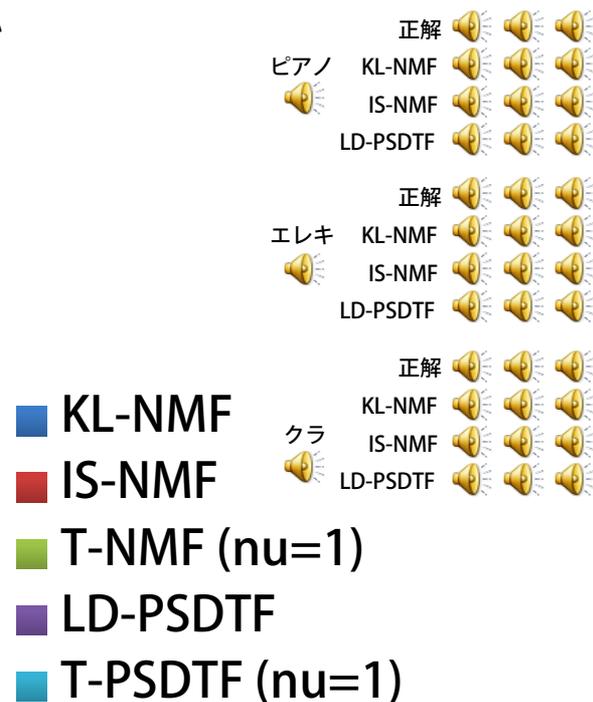
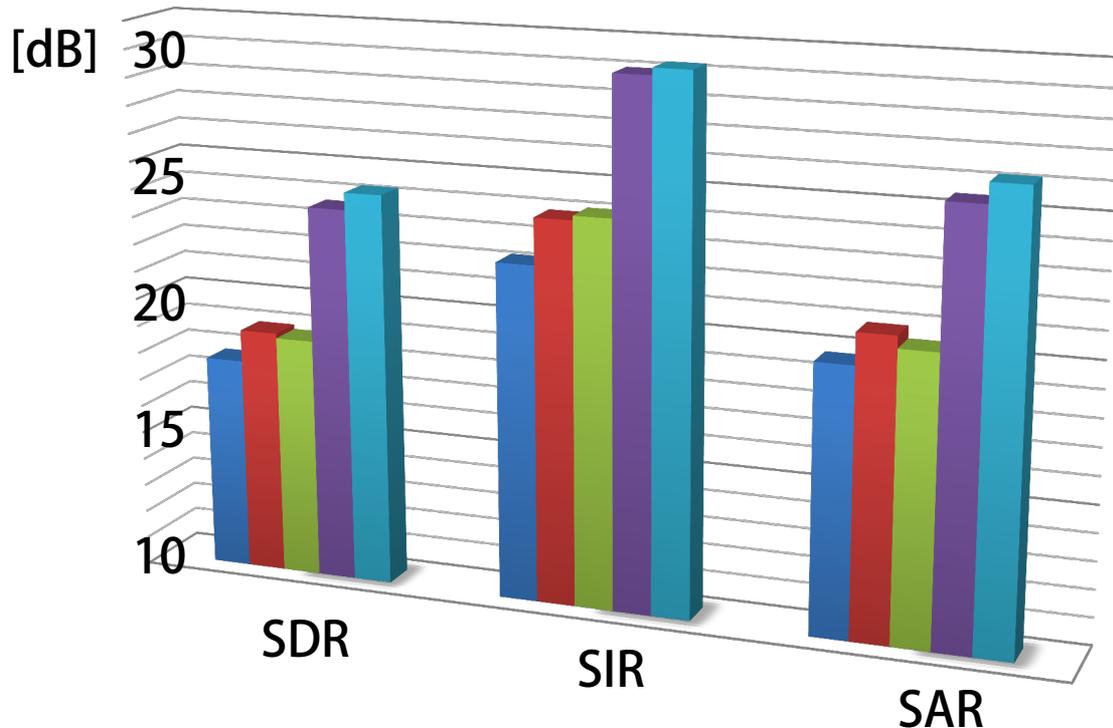
# 評価実験

- 人工的に混合音を作成して基本的な音源分離性能を評価
  - 基底数は $K=3$  (C4, E4, G4に対応)
  - RWC研究用音楽データベースから3種類の楽器の単独音を使用
    - ピアノ・エレキギター (減衰音) / クラリネット (持続音)
  - BSS Eval Toolbox [Vincent2006] を使用して評価



# 音源分離結果

- $t$ -PSDTFは優れた音源分離品質を達成・さらなる調査が必要
  - IS-NMFより $t$ -NMFの方がわずかに精度が低いが局所解に陥りにくい傾向
  - LD-PSDTFより $t$ -PSDTFの方がわずかに精度が高い



# まとめと今後の課題

- Done:  $t$ 過程に基づく半正定値テンソル分解 (PSDTF)
  - ガウス過程( $\nu = \infty$ )とコーシー過程( $\nu = 1$ )を特殊形を含む
    - 再生性が成立するのはこれらのみ
  - 過学習抑制効果で局所解を回避しやすい傾向
    - $t$ -NMFはIS-NMFの局所解問題の解決策？
    - $\nu$ を徐々に上昇させるアニーリングが有効？
- Todo:  $S\alpha S$ 過程に基づく半正定値テンソル分解
  - ガウス過程とコーシー過程を特殊形を含む
    - 再生性が全ての $\alpha$ で成立
  - より音響信号にフィットしたモデル化が可能

