

ロボットのための マイクアレイ音響信号処理

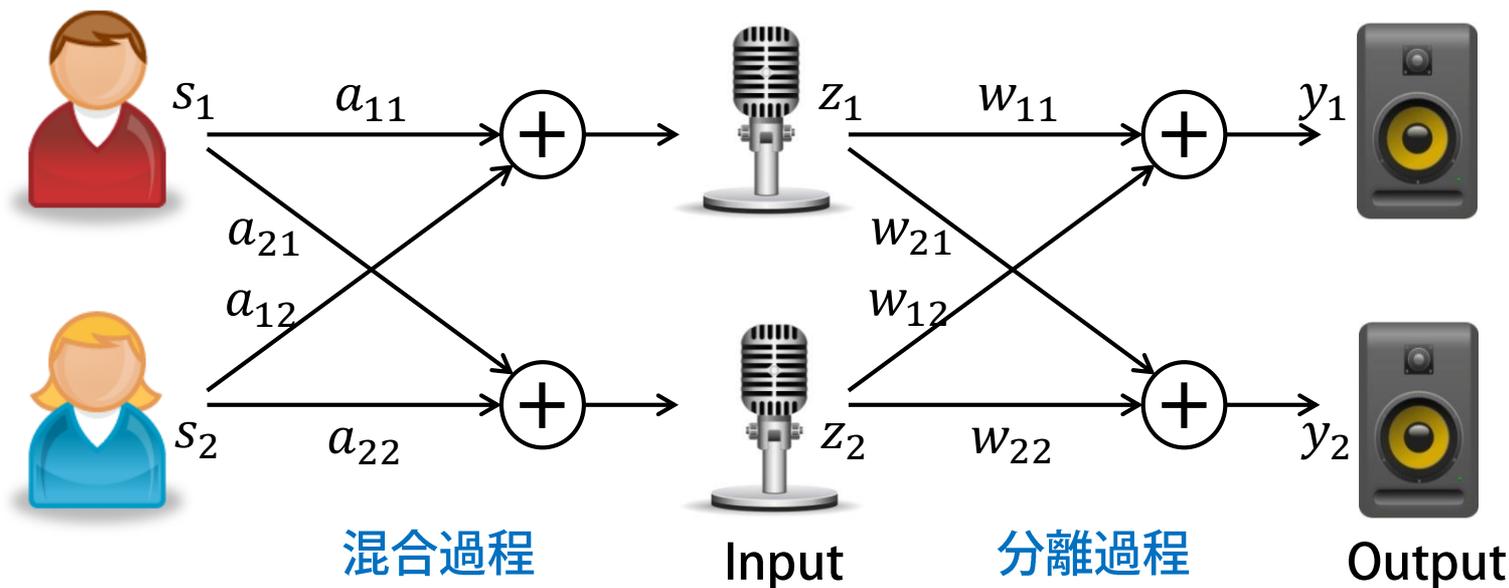
京都大学 大学院情報学研究科 知能情報学専攻
理研 AIP 音響情景理解チーム (兼任)

吉井 和佳

yoshii@kuis.kyoto-u.ca.jp

- マルチチャネルの混合音に対する音源分離と音源定

- 入力： z_1, z_2, \dots, z_N 出力： y_1, y_2, \dots, y_M ($\approx s_1, s_2, \dots, s_M$)
 - 混合過程： $s_1, s_2, \dots, s_M \rightarrow$ 観測 z_1, z_2, \dots, z_N
 - 条件： A が既知 $\leftrightarrow A$ が未知 (ブラインド)



- ブラインド条件かどうかで利用できる技術が異なる
 - 非ブラインド条件
 - ◆ ビームフォーミング
 - ◆ MUSIC (multiple signal classification) 法
 - ブラインド条件
 - ◆ 独立成分分析 (ICA) ・ 独立ベクトル分析 (IVA)
 - ◆ マルチチャネル非負値行列分解 (NMF)
 - ◆ 時間周波数クラスタリング ・ マスキング
 - より先進的なトピック
 - ◆ 統一的な確率モデルによる分離 ・ 定位 ・ 残響除去
 - ◆ ノンパラベイズモデルによる音源数の自動推定

音源信号と観測信号の関係

- 各チャネルの観測信号は音源信号の時間遅れ
 - 音源信号 $s(t)$ を M 個のマイクで観測するとする
 - 各マイク m では遅れ τ_m が発生する

$$\mathbf{z}(t) = \begin{bmatrix} z_1(t) \\ z_2(t) \\ \vdots \\ z_M(t) \end{bmatrix} = \begin{bmatrix} s(t - \tau_1) \\ s(t - \tau_2) \\ \vdots \\ s(t - \tau_M) \end{bmatrix} \xrightarrow{\text{フーリエ変換}} \mathbf{z}(\omega) = \begin{bmatrix} Z_1(\omega) \\ Z_2(\omega) \\ \vdots \\ Z_M(\omega) \end{bmatrix}$$

観測信号

$$Z_m(\omega) \equiv \int_{-\infty}^{\infty} z_m(t) e^{-j\omega t} dt = \int_{-\infty}^{\infty} s(t - \tau_m) e^{-j\omega t} dt = e^{-j\omega \tau_m} S(\omega)$$

$$S(\omega) \equiv \int_{-\infty}^{\infty} s(t) e^{-j\omega t} dt \quad \mathbf{a}(\omega) = \begin{bmatrix} a_1 \\ \vdots \\ a_M \end{bmatrix} \equiv \begin{bmatrix} e^{-j\omega \tau_1} \\ \vdots \\ e^{-j\omega \tau_M} \end{bmatrix}$$

ステアリングベクトル

$$\mathbf{z}(\omega) = \mathbf{a}(\omega) S(\omega)$$

観測

関係

音源

- マイク位置

$$\mathbf{p}_m = \left[\left((m-1) - \frac{M-1}{2} \right) d_x, 0, 0 \right]^T$$

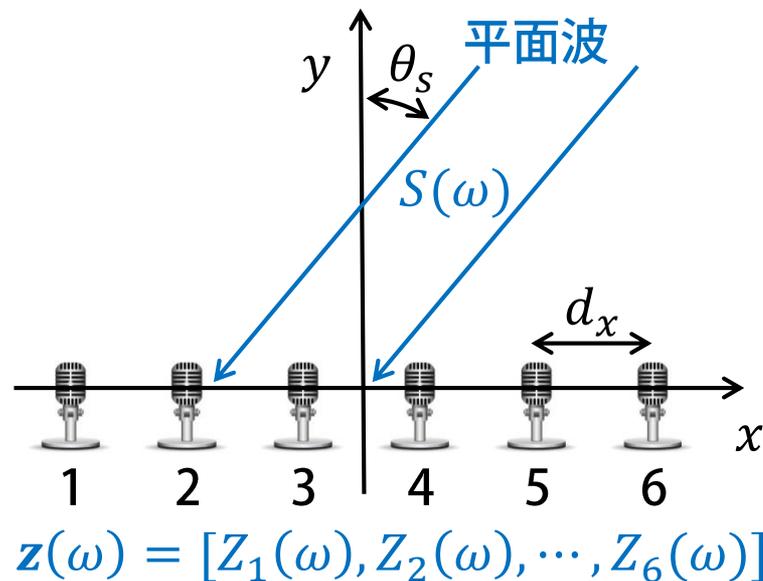
- 遅延

$$\tau_m = - \left((m-1) - \frac{M-1}{2} \right) \frac{d_x}{c} \sin \theta_s$$

- ステアリングベクトル

$$\mathbf{a}_m(\omega) = \exp \left(j \left((m-1) - \frac{M-1}{2} \right) \frac{2\pi d_x}{\lambda} \sin \theta_s \right)$$

$$\mathbf{a}(\omega) = e^{-\frac{j(M-1)\psi}{2}} [1, e^{j\psi}, e^{j2\psi}, \dots, e^{j(M-1)\psi}]^T \quad \left(\psi = \frac{2\pi d_x}{\lambda} \sin \theta_s \right)$$



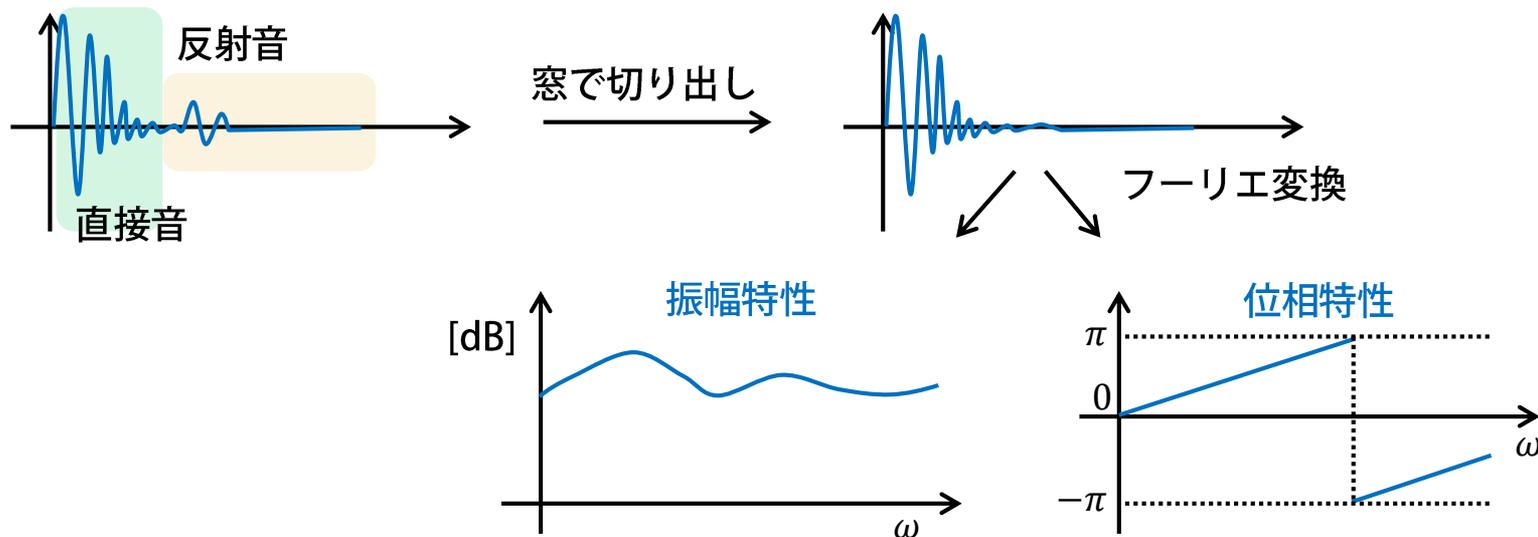
$$\mathbf{z}(\omega) = \mathbf{a}(\omega) S(\omega)$$

ステアリングベクトルの準備

- 幾何的に計算する (中空を仮定する場合が多い)
 - 公式使うだけ: $a_m(\omega) = \exp(-jk^T \mathbf{p}) = \exp\left(j \frac{2\pi}{\lambda} \mathbf{u}^T \mathbf{p}_m\right)$
- 実際に計測してみる
 - 直接音のみを抽出する
 - フーリエ変換を行って周波数軸に変換しておく

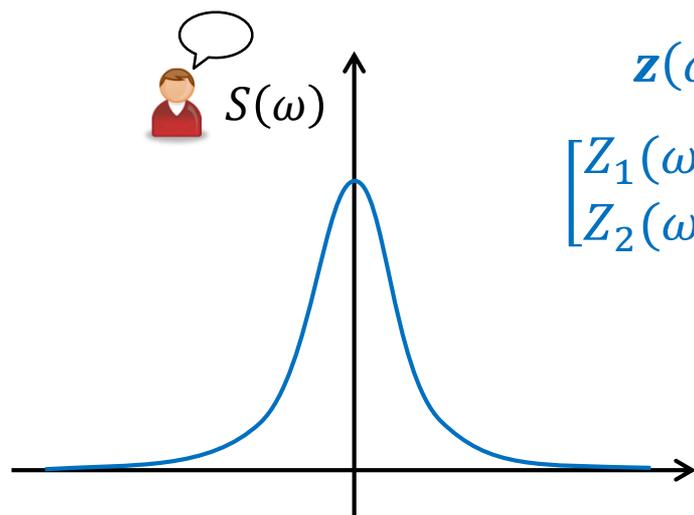
音源位置

マイクの位置



単一音源の観測と確率モデル

- 確率モデル $y = f(x)$ をきめて x の分布を与える
 - 観測可能な y が従う分布が分かる
 - 必要に応じてヤコビアンを計算する



We often assume
 $S(\omega) \sim N_c(0,1)$

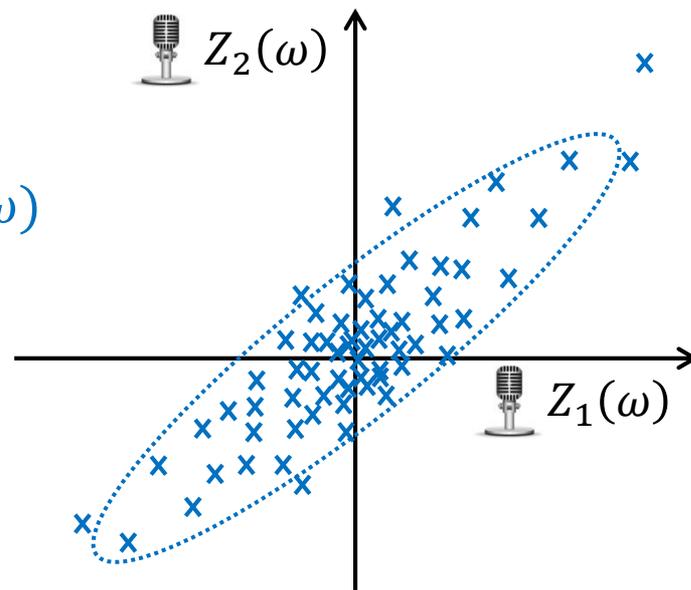
$$\mathbf{z}(\omega) = \mathbf{a}(\omega)S(\omega)$$

$$\begin{bmatrix} Z_1(\omega) \\ Z_2(\omega) \end{bmatrix} = \begin{bmatrix} a_1(\omega) \\ a_2(\omega) \end{bmatrix} S(\omega)$$



$$\mathbf{a}(\omega) = \begin{bmatrix} a_1 \\ \vdots \\ a_M \end{bmatrix} \equiv \begin{bmatrix} e^{-j\omega\tau_1} \\ \vdots \\ e^{-j\omega\tau_M} \end{bmatrix}$$

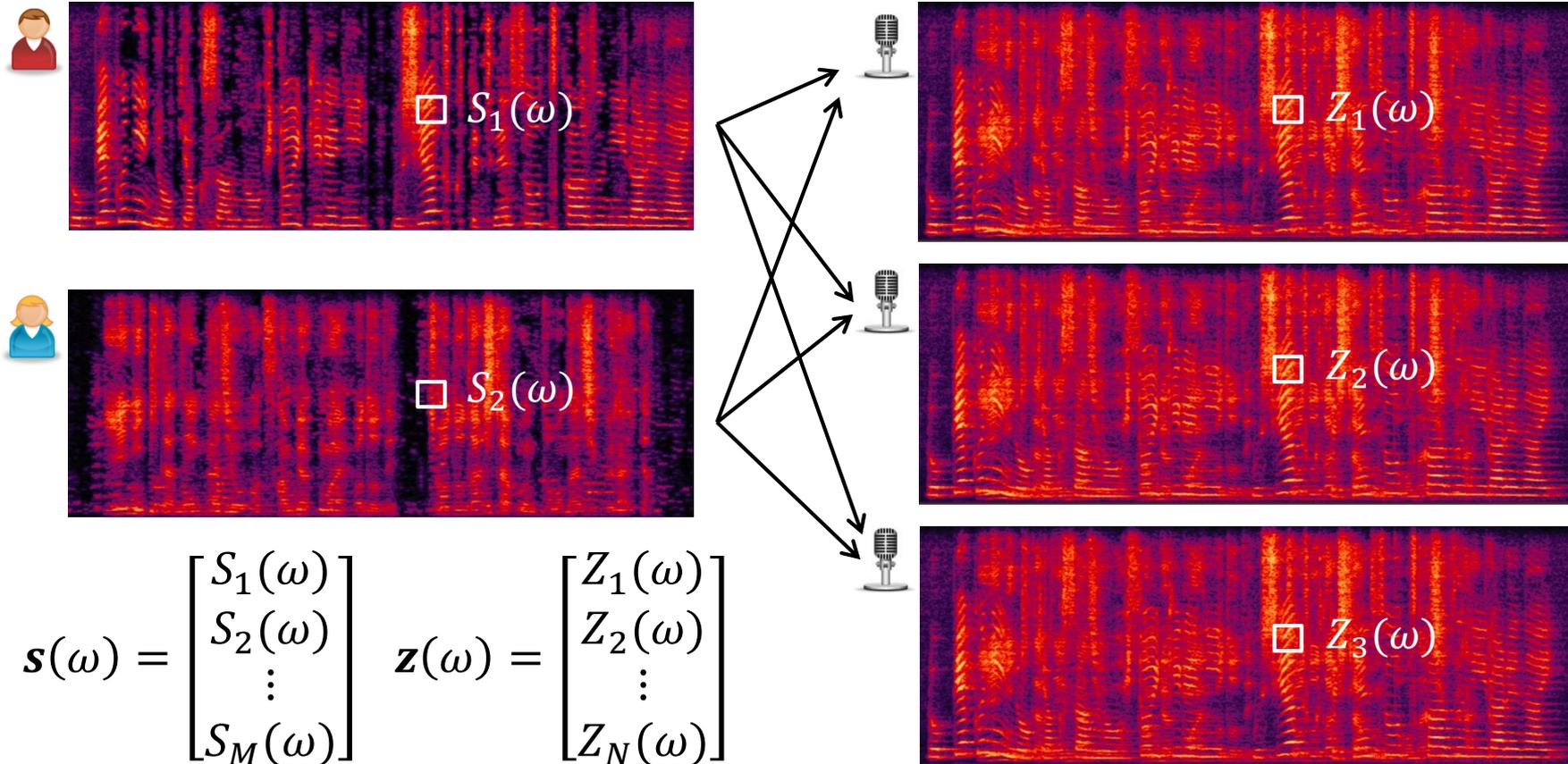
ステアリングベクトル



$\mathbf{z}(\omega) \sim N_c(0, \mathbf{a}(\omega)\mathbf{a}^H(\omega))$
空間相関行列

複数マイクによる混合音の観測

- N 個の音源を M 個のマイクで観測する



- 観測信号は N 個の音源“イメージ”の線形和
 - 音源信号はモノラルでも、複数マイクで観測するとマルチチャンネル
 - マイクアレイの地点における音源のイメージとよぶ

単一音源

複数音源

$$\mathbf{z}_S(\omega) = \mathbf{a}(\omega)S(\omega) \quad \longrightarrow \quad \mathbf{z}_S(\omega) = \sum_{i=1}^N \mathbf{a}_i(\omega)S_i(\omega) = \mathbf{A}(\omega)\mathbf{s}(\omega)$$

$$\mathbf{z}_S(\omega) = \begin{bmatrix} Z_{S1}(\omega) \\ Z_{S2}(\omega) \\ \vdots \\ Z_{SM}(\omega) \end{bmatrix}$$

$M \times N$ の混合行列

$$\mathbf{A}(\omega) = [\mathbf{a}_1(\omega), \dots, \mathbf{a}_N(\omega)]$$

$$\mathbf{s}(\omega) = \begin{bmatrix} S_1(\omega) \\ S_2(\omega) \\ \vdots \\ S_N(\omega) \end{bmatrix}$$

$\mathbf{a}_n(\omega)$: 音源 n に対する
ステアリングベクトル

- 物理的現象 $z = As + v$ に対する確率モデルをつくりたい

- 決定的な信号モデル

$$p(v) = N(v|\mathbf{0}, K) \xrightarrow[\substack{\text{線形変換} \\ z = As + v}]{\quad} \text{尤度関数} : p(z; \Theta) = N(z|As, K)$$

Find $\Theta = \{A, s, K\}$ that maximizes $p(z; \Theta)$

- 確率的な信号モデル

A は N 個の音源の方向 $\{\theta_1, \dots, \theta_N\}$ に依存
 Γ は N 個の音源のパワー $\{\gamma_1, \dots, \gamma_N\}$ に依存

$$p(v) = N(v|\mathbf{0}, K) \xrightarrow[\substack{\text{線形変換} \\ z = As + v}]{\quad} \text{尤度関数} : p(z; \Theta) = N(z|\mathbf{0}, A\Gamma A^H + K)$$
$$p(s) = N(s|\mathbf{0}, \Gamma)$$

Find $\Theta = \{A, \Gamma, K\}$ that maximizes $p(z; \Theta)$

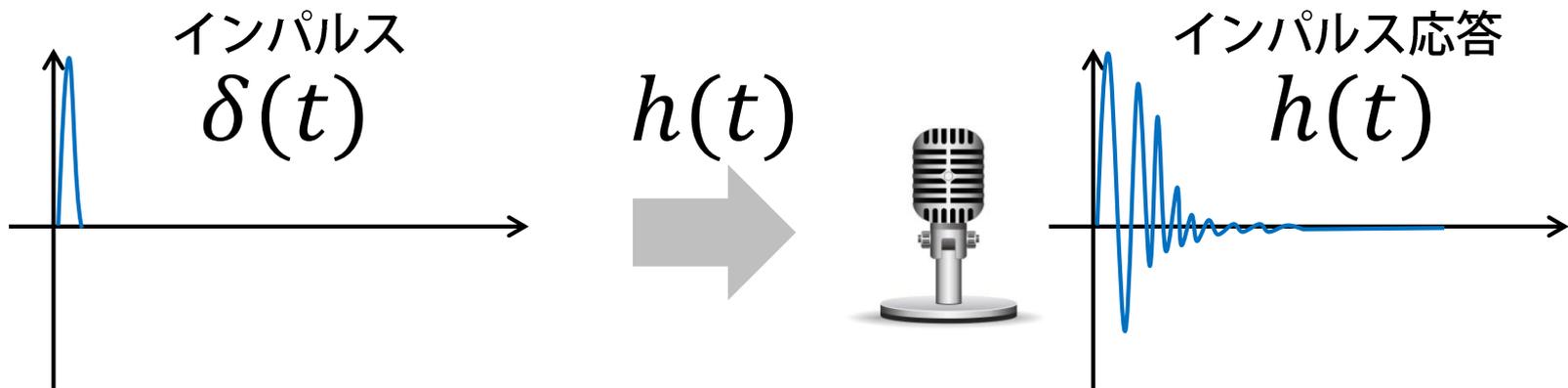
$$\Gamma = E[ss^H] (= \text{diag}(\gamma_1, \dots, \gamma_N))$$

事前分布 $p(\Theta)$ を導入すれば
ベイズ的な取り扱いも可能

$$p(\Theta|z) = \frac{p(z|\Theta)p(\Theta)}{p(z)}$$

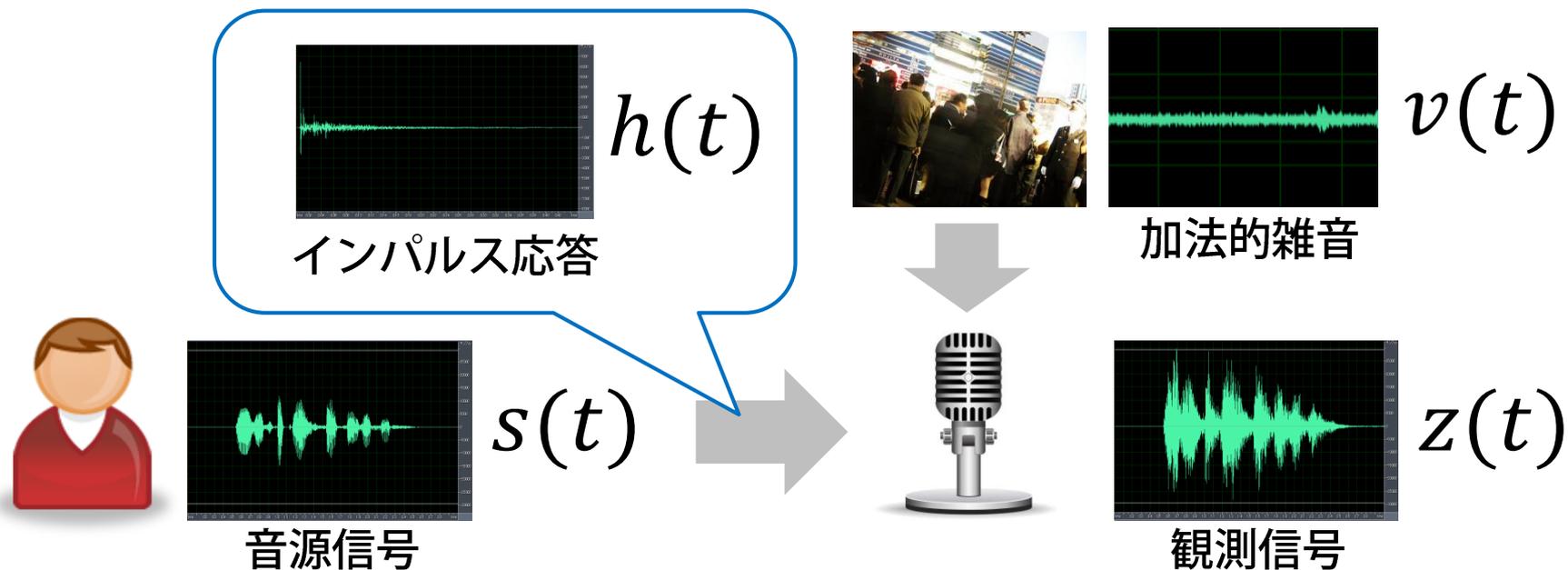
インパルス応答

- 音源からインパルス信号が生成された場合に、各マイクで観測される信号のこと
 - 音源信号が伝播するまで、反射やまわりこみなど様々な影響を受ける
 - インパルス応答 (時間領域) = 周波数応答/伝達関数 (周波数/z領域)
 - 部屋の音響特性を表しているので、シミュレーションに有用



$$h(t) = h(t) * \delta(t)$$

- 部屋の音響特性はしばしば線形モデルで記述される
 - 音源信号 + 音響特性 + 加法的雑音 → 観測信号



$$z(t) = h(t) * s(t) + v(t)$$

部屋の音響特性のシミュレーション

- インパルス応答を用いれば、音源信号を与えると、その部屋で観測されるであろう信号をシミュレートできる

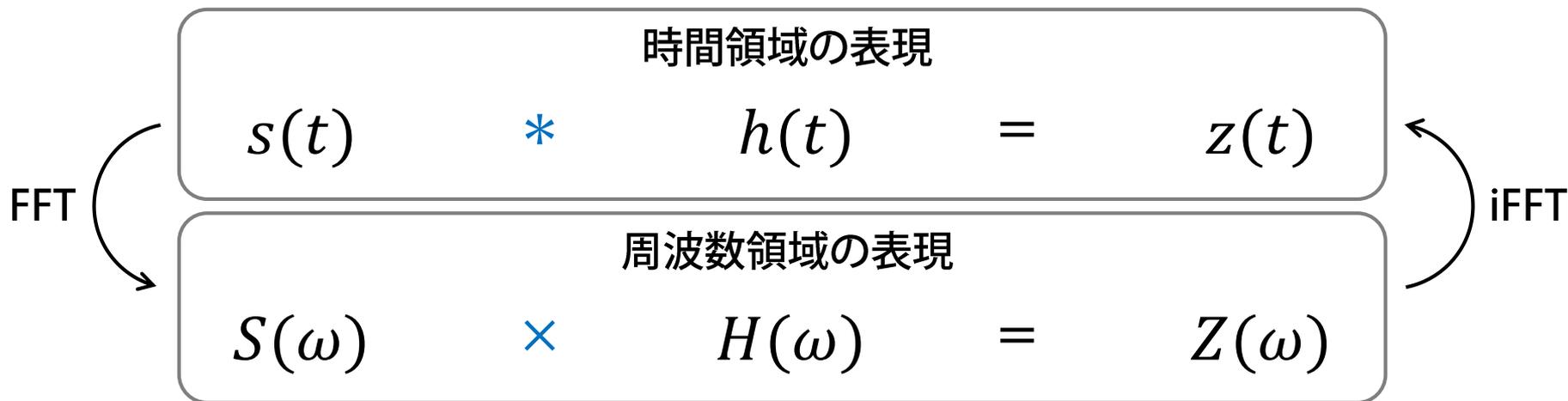
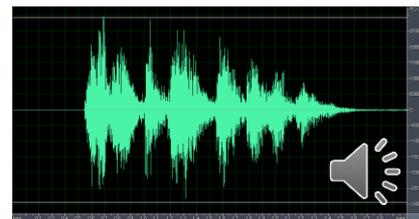
音源信号



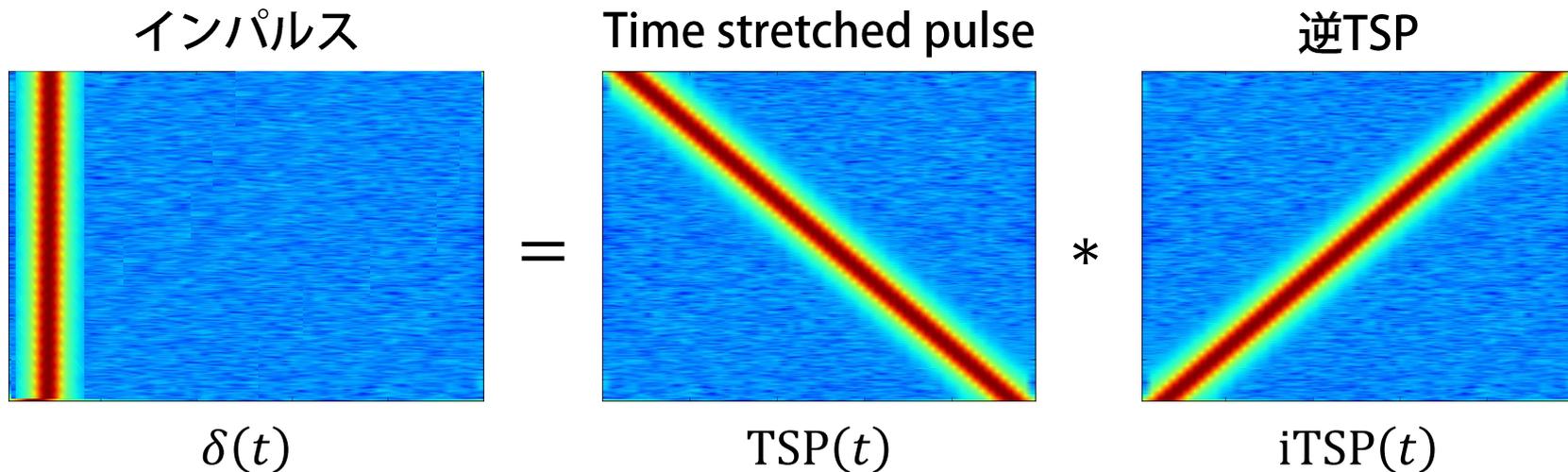
インパルス応答



音響信号

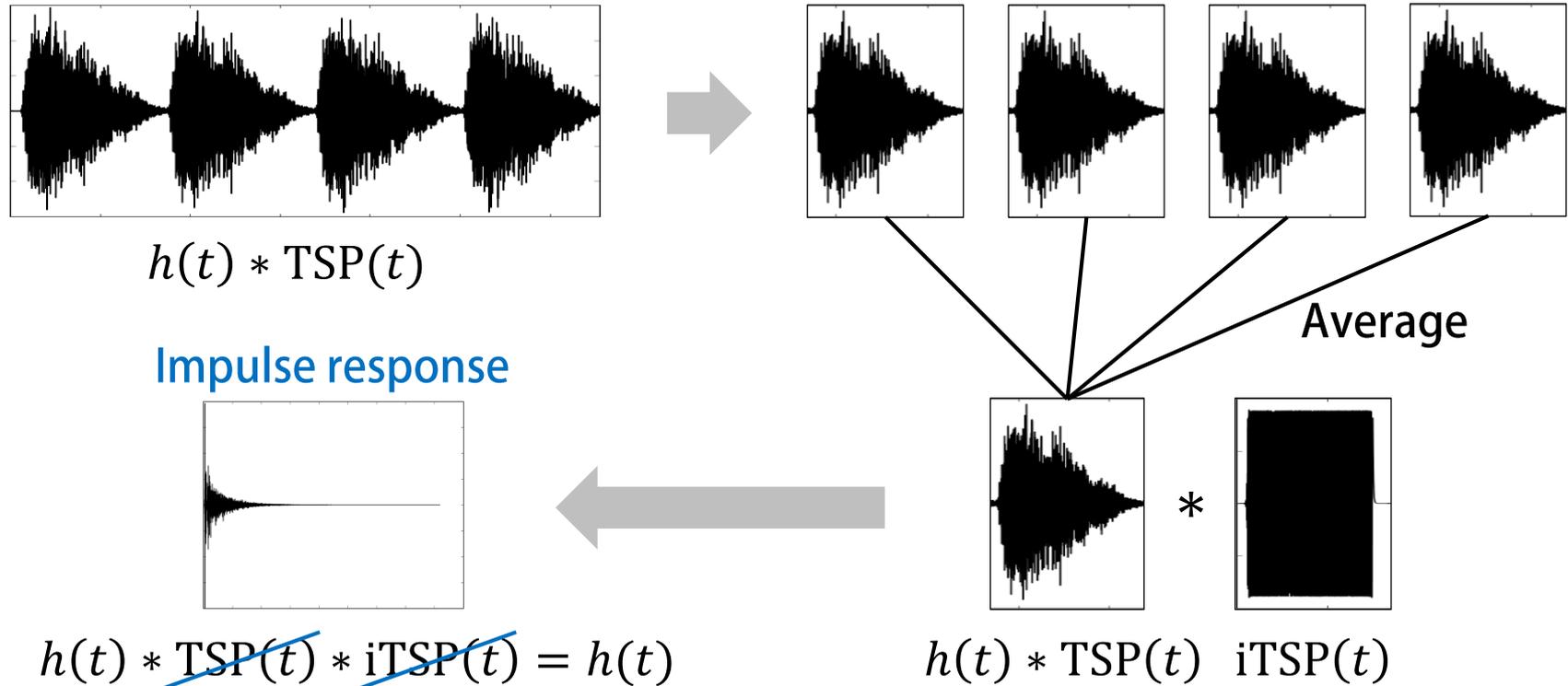


- インパルス応答を測定するにはTSP信号を利用する
 - 安直な方法：インパルス信号を出力して、インパルス応答を測定
 - 全周波数帯域に等しいパワーがあるので再生が困難
 - 工夫した方法：TSP応答を測定してからインパルス応答に変換する
 - TSP信号は再生しやすく、逆TSPの畳み込みも容易

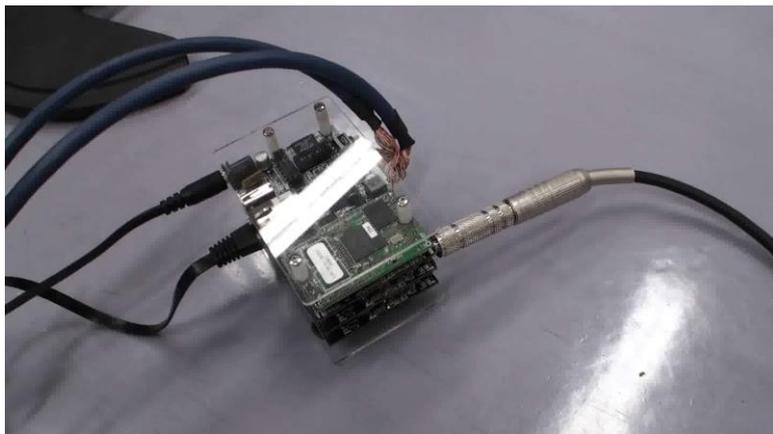
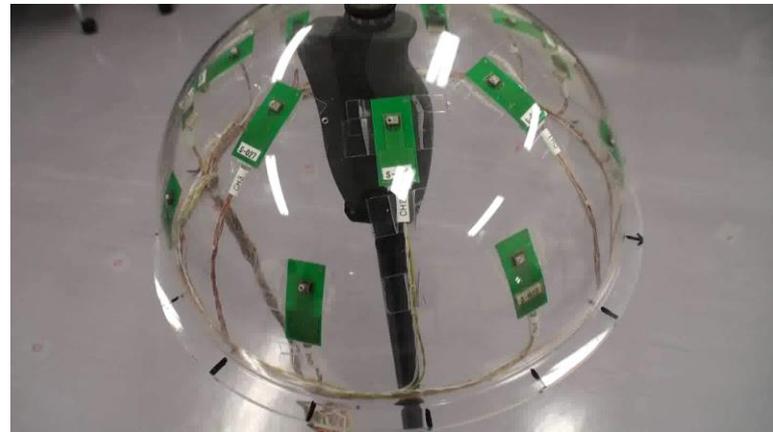


インパルス応答の測定

- TSP応答を何度か測定・平均してから、逆TSPを畳み込み
 - TSPとiTSPがキャンセルされてインパルス応答がでてくる



- TSP信号の生成と録音を行うための機材を準備する



スピーカーと耳栓

TSP信号でも音量がかなり大きい

マイクアレイ

各マイクでTSP応答を録音する

録音機材

すべてのマイクを同期させておく

TSP Recording



独立成分分析

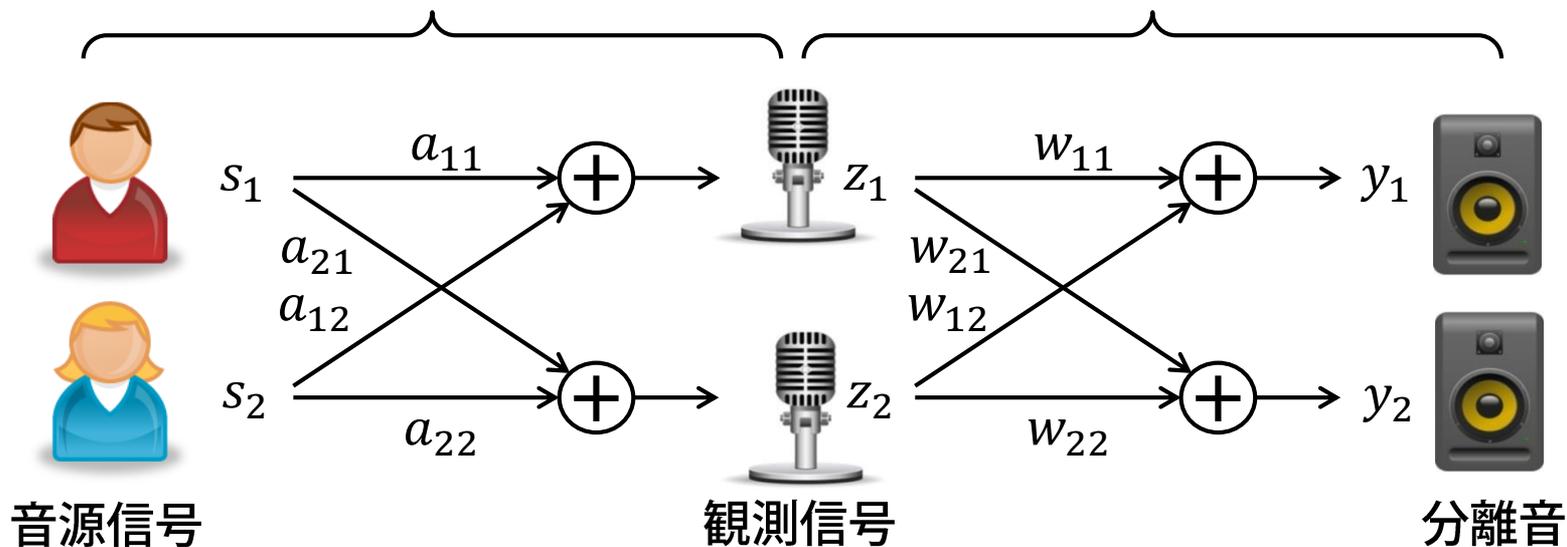
- 周波数領域において瞬時混合過程をモデル化する

- N 個の音源を M 個のマイクで観測する $\{M = N$ を仮定

$$\mathbf{z} = \mathbf{A}\mathbf{s} = \sum_{i=1}^N \mathbf{a}_i s_i \quad \mathbf{y} = \mathbf{W}\mathbf{z} = \mathbf{W}\mathbf{A}\mathbf{s} \quad \rightarrow \quad \text{if } \mathbf{W} = \mathbf{A}^{-1}, \mathbf{y} \approx \mathbf{s}$$

混合過程 : $\mathbf{z} = \mathbf{A}\mathbf{s}$

分離過程 : $\mathbf{y} = \mathbf{W}\mathbf{z}$



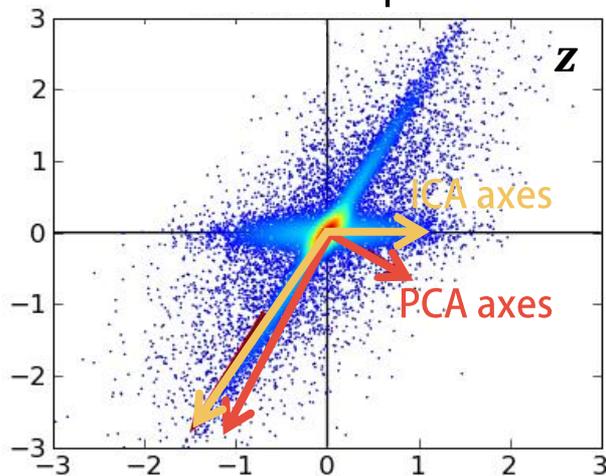
主成分分析 (PCA) と独立成分分析 (ICA)

- PCAは互いに無相関となるような次元に変換
 - 軸が直交していればOK
- ICAは互いに独立となるような次元に変換
 - 軸同士が独立でないといけない

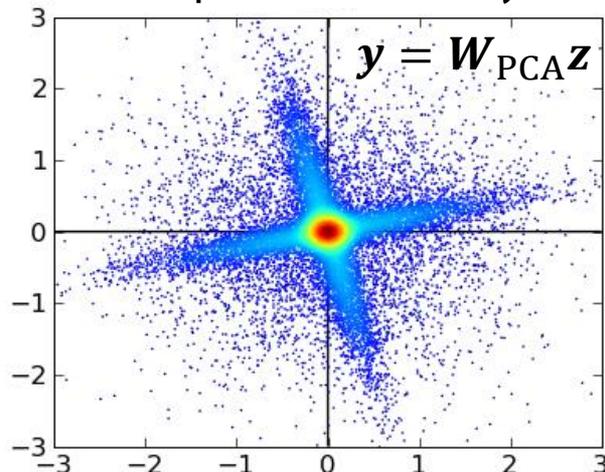
← 十分条件

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix} \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} = \mathbf{w}_1 z_1 + \mathbf{w}_2 z_2$$

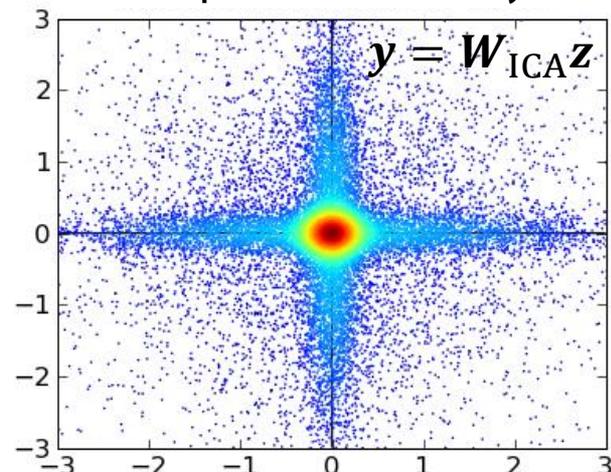
Observed space



Latent space discovered by PCA



Latent space discovered by ICA

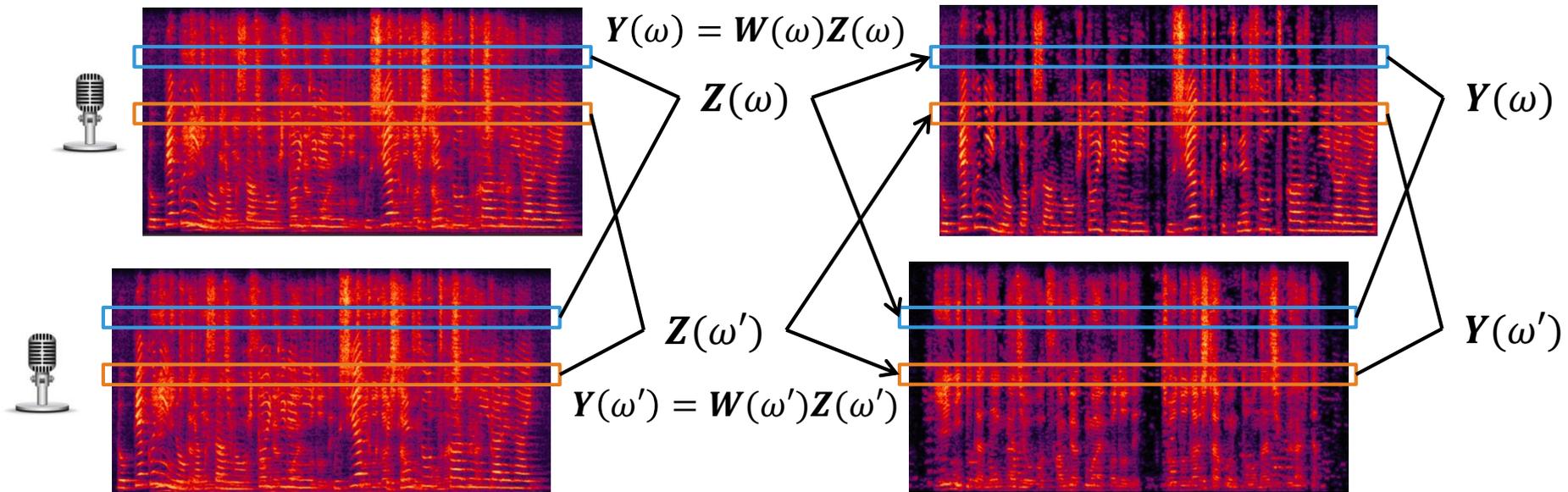


- パーミュテーション問題

- 各周波数で別々にICAを行うので、音源ごとにまとめる必要がある

- 振幅の曖昧性の問題

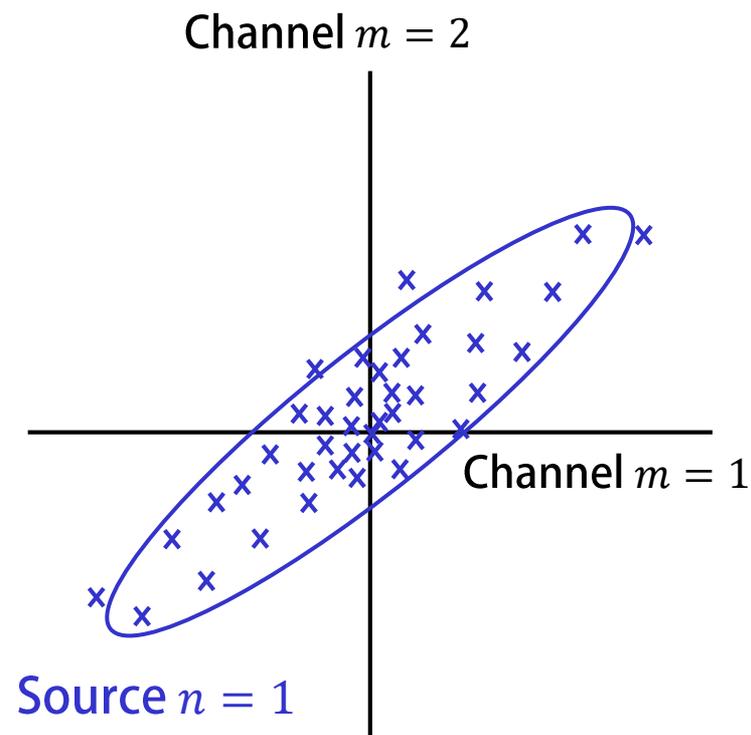
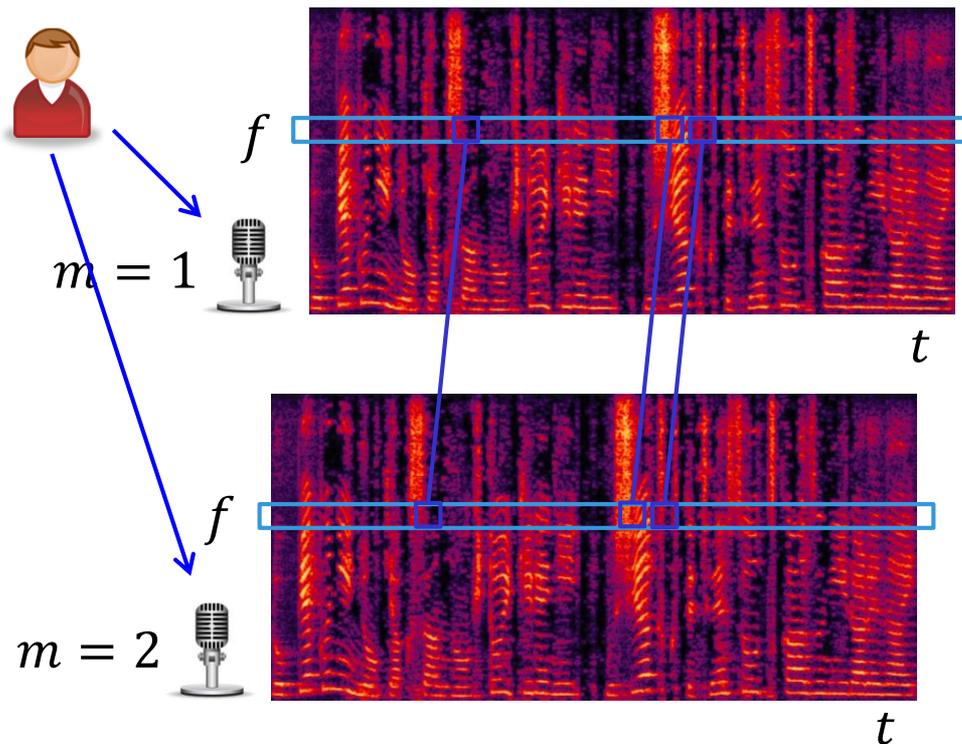
- 各周波数ごとの振幅のスケールもそろえる必要がある



時間周波数クラスタリング

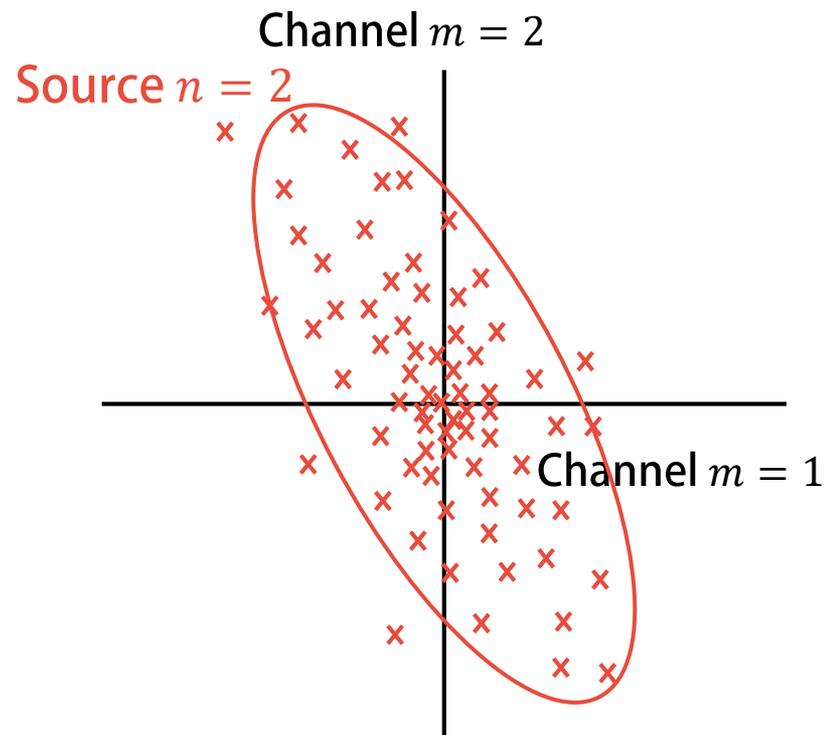
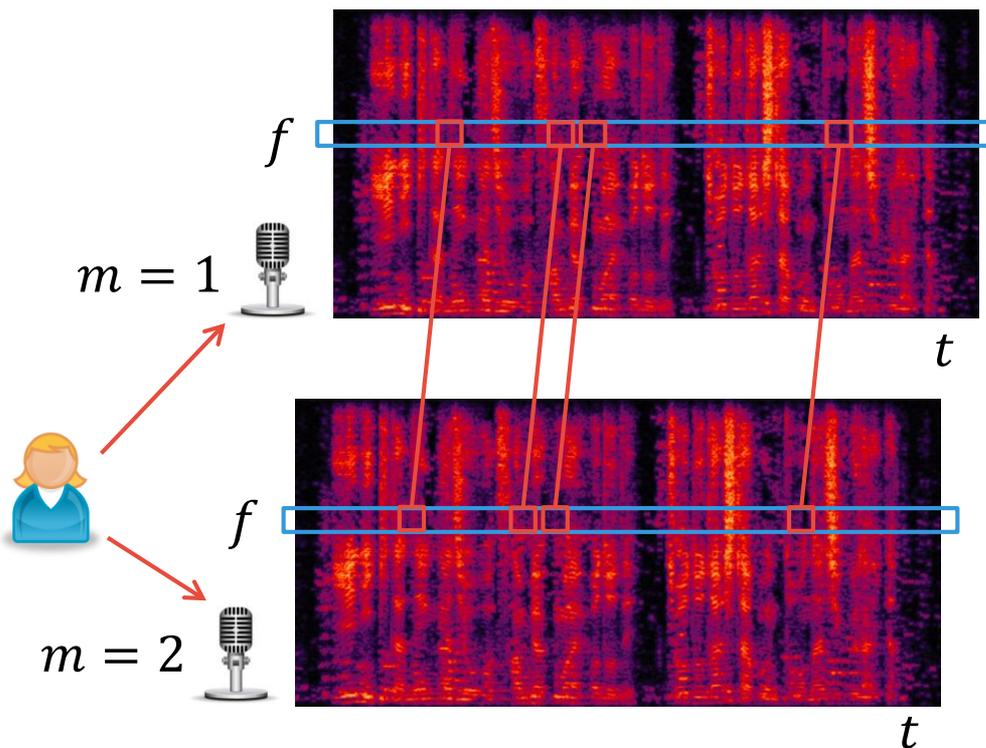
- 各音源由来のイメージは明確な空間的な構造を持つ
 - 空間相関行列をもつガウス分布に従う

観測データ： $x_{tf} = [x_{tf1}, x_{tf2}, \dots, x_{tfM}]$



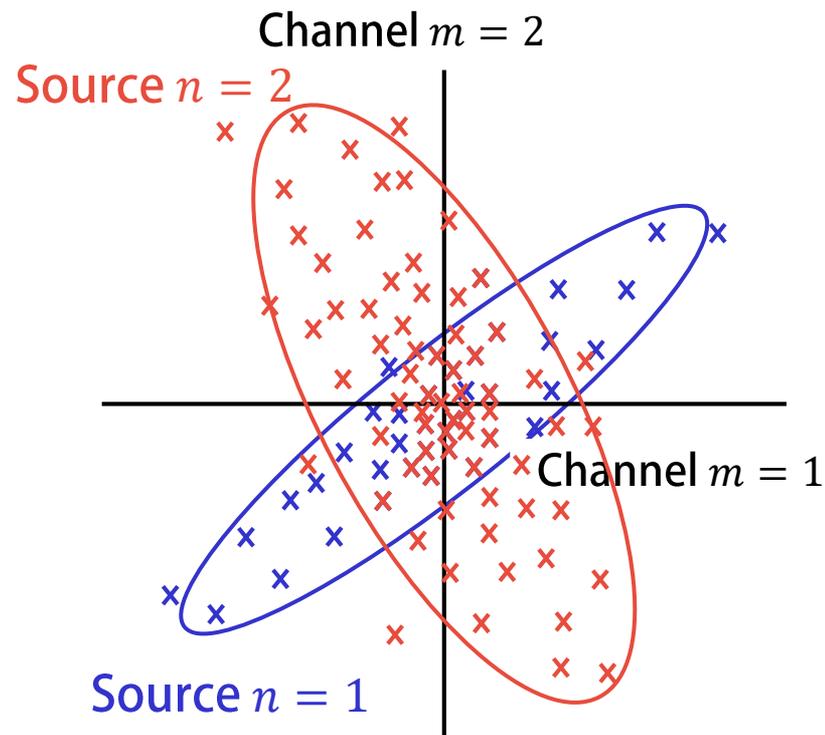
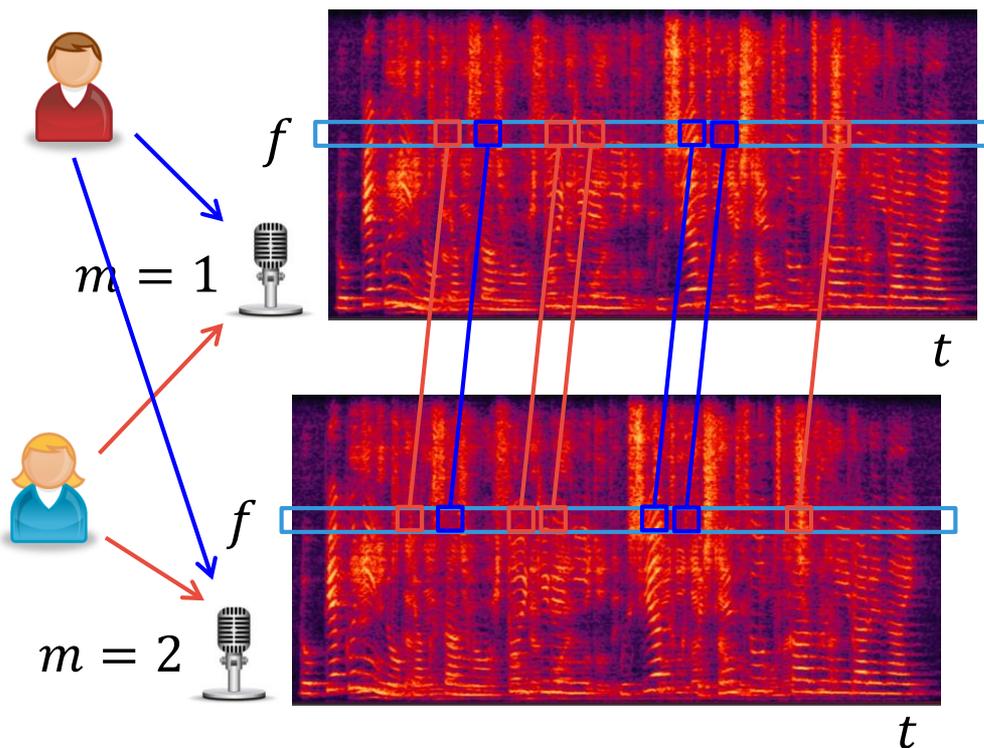
- 各音源由来のイメージは明確な空間的な構造を持つ
 - 空間相関行列をもつガウス分布に従う

観測データ： $x_{tf} = [x_{tf1}, x_{tf2}, \dots, x_{tfM}]$

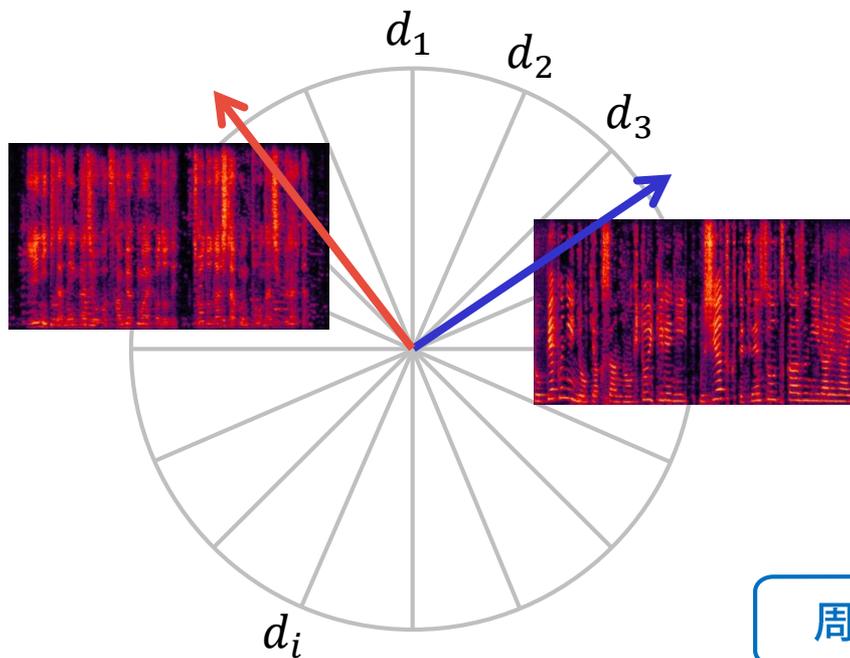


- 観測スペクトルは複数の空間的な構造がまじっている
 - 各ビンにおいてはいずれか一つの音源が優勢と仮定してみる

Observed data: $x_{tf} = [x_{tf1}, x_{tf2}, \dots, x_{tfM}]$



- 各時間・周波数ビンをいずれかの音源にクラスタリング
 - 時間 t ・周波数 f が音源 k に属する： $z_{tf} = k$
 - \mathbf{H}_{fd} ：周波数 f ・方向 d に関する空間相関行列



観測モデル [Duong 2010]

$$\mathbf{x}_{tf} \sim N_c \left(\mathbf{x}_{tf} \mid \mathbf{0}, \left(\lambda_{tf} \mathbf{H}_{fd_{z_{tf}}} \right)^{-1} \right)$$

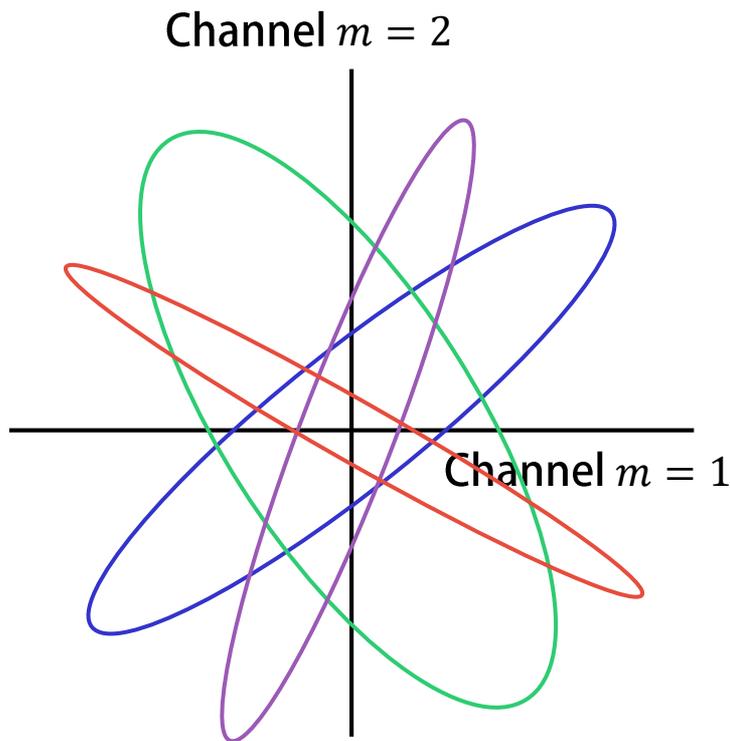
時間 t ・周波数 f における優勢な音源

事前分布の導入とベイズ推論 [Otsuka 2014]

$$\mathbf{H}_{fd} \sim W_c \left(\left(\mathbf{a}_{fd} \mathbf{a}_{fd}^H + \epsilon \mathbf{I} \right)^{-1}, \nu_0 \right)$$

周波数 f ・方向 d におけるステアリングベクトル

- 観測信号に合わせて適切な個数の音源信号を推定
 - 理論的には無限個存在するが、そのうち一部だけが実体化される



観測モデル [Duong 2010]

$$\mathbf{x}_{tf} \sim N_c \left(\mathbf{x}_{tf} \mid \mathbf{0}, \left(\lambda_{tf} \mathbf{H}_f d_{z_{tf}} \right)^{-1} \right)$$

時間 t ・ 周波数 f における優勢な音源

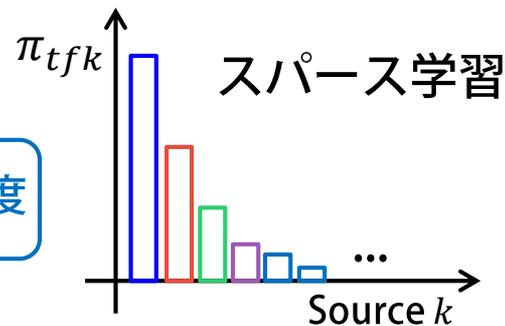
階層ディリクレ過程 ($k \rightarrow \infty$) [Otsuka 2014]

$$\boldsymbol{\pi}_{tf} \sim \text{HDP}(\alpha, \gamma, \boldsymbol{\beta})$$

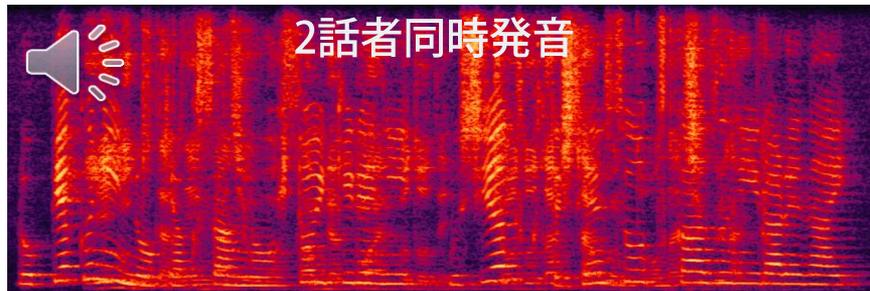
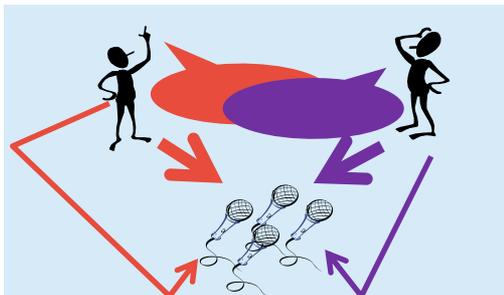
集中度

基底速度

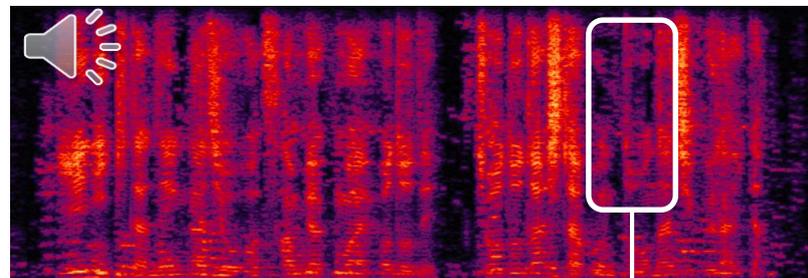
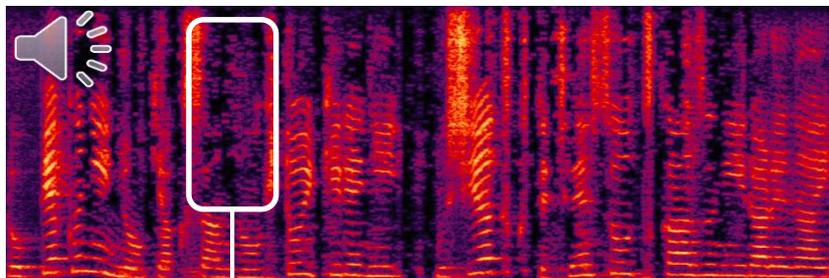
$$z_{tf} \sim \text{Categorical}(\boldsymbol{\pi}_{tf})$$



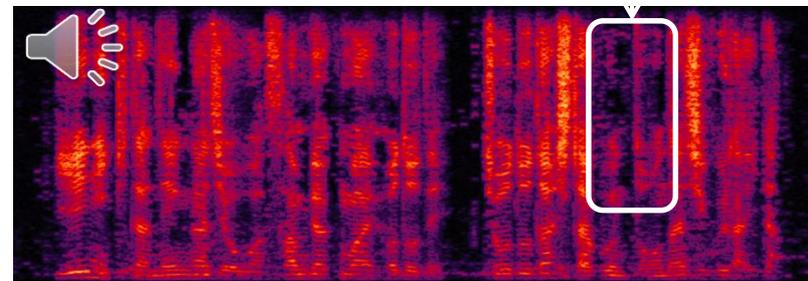
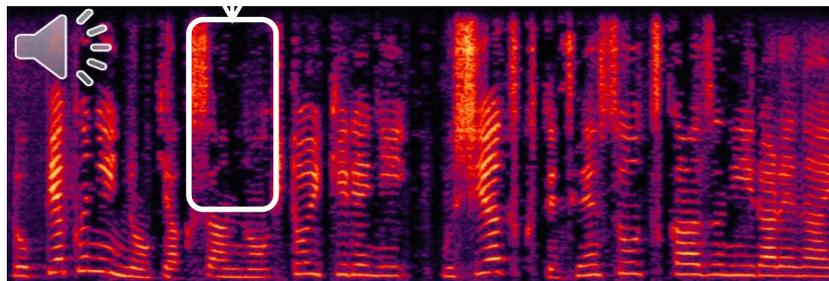
定位 + 分離 + 残響除去



残響除去なし



残響除去あり



京大時計台国際交流ホール



マイクロホンアレイ



観測した混合音



分離音



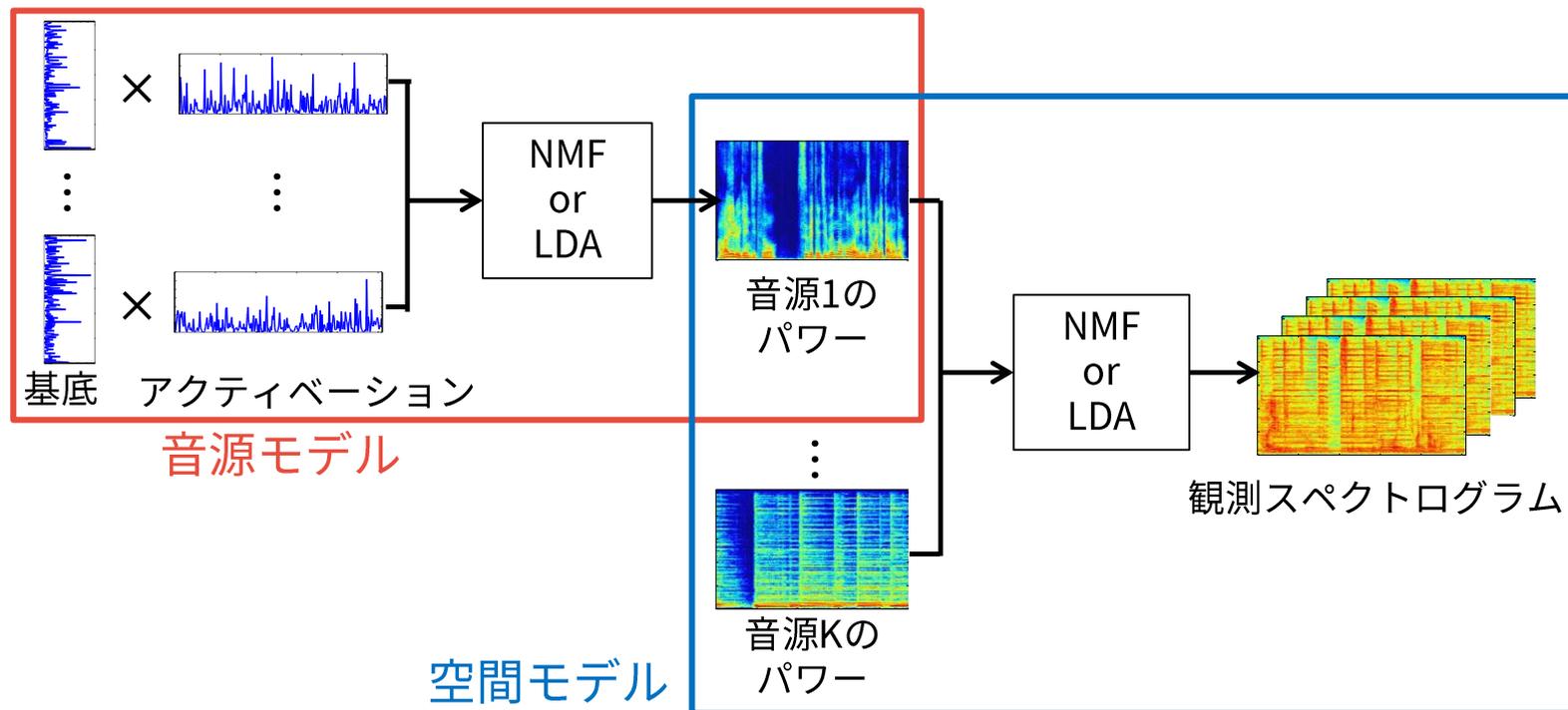
周囲の人の話し声

背景雑音

音源信号と観測信号の 階層ベイズモデリング

観測スペクトログラムの階層的生成モデル

- 音源モデル：音源信号のパワーの生成過程をモデル化
- 空間モデル：音源イメージの位相をモデル化



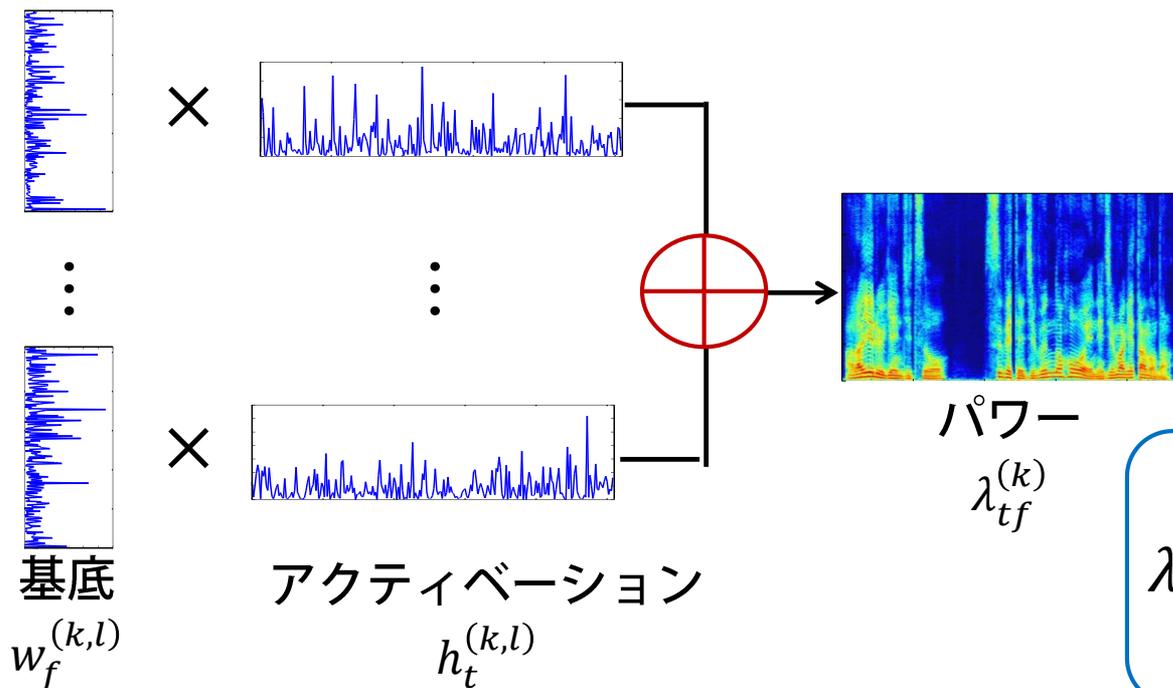
- 音源信号の生成過程と音源信号の重畳過程を同時に考慮
 - 音源モデル：各音源の音源スペクトログラムの構造を表現
 - ◆ 低ランク性が有用 (e.g., NMF, LDA)
 - 空間モデル：複数の音源がどのように混じるかを表現
 - ◆ 従来の空間相関行列が手がかかり

空間モデル

	因子モデル (factor)	混合モデル (mixture)
因子モデル(factor)	factor-factor	factor-mixture
混合モデル (mixture)	mixture-factor	mixture-mixture

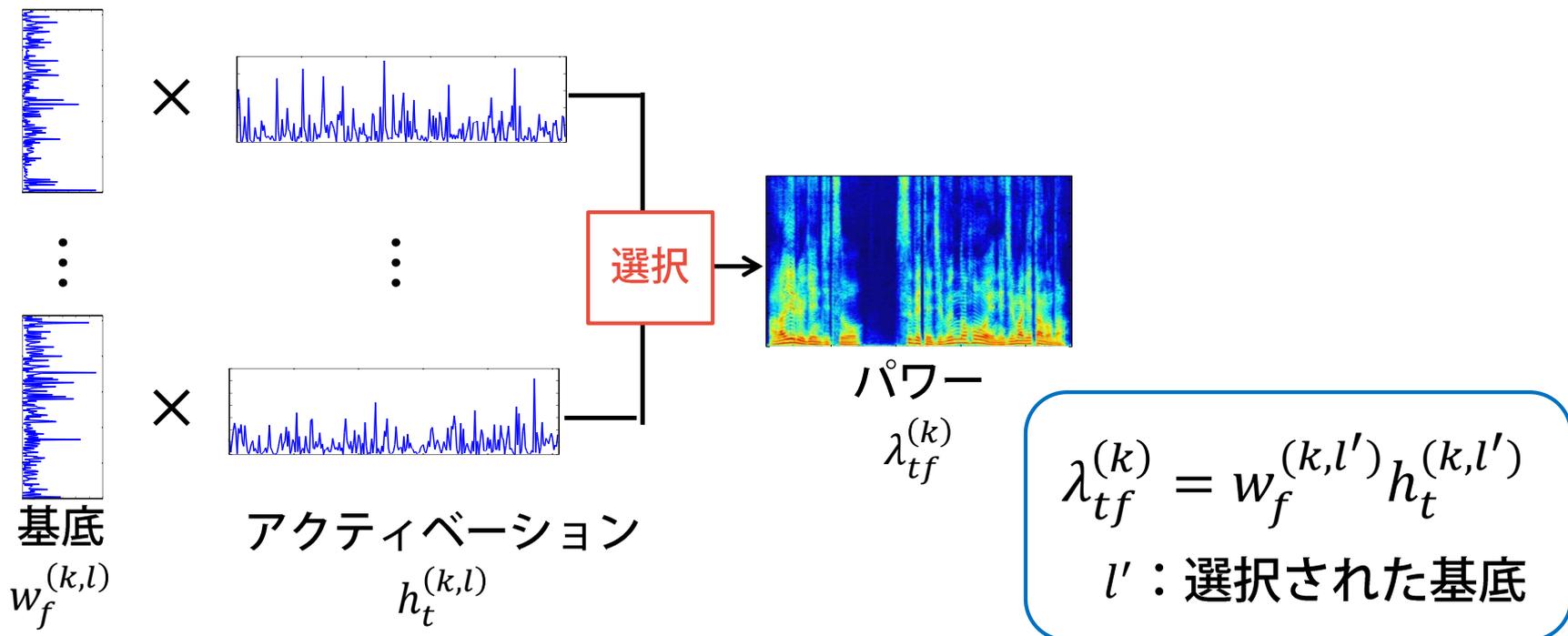
音源モデル：因子モデル (NMF)

- 音源スペクトログラムを基底スペクトログラムの足し合わせで表現
 - 音源 k の時間 t ・周波数 f におけるパワー λ_{tf}^k を、基底とアクティベーションの積のすべてを使って表現

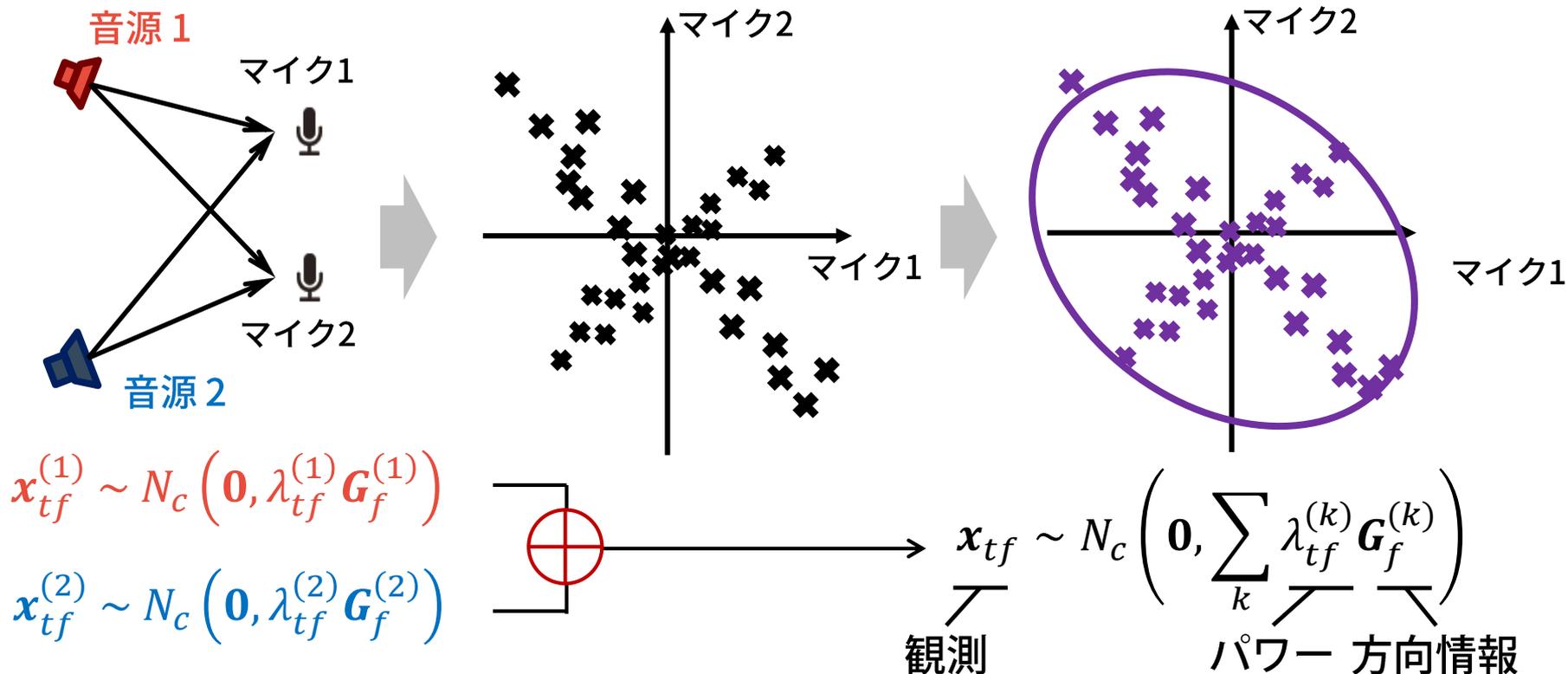


$$\lambda_{tf}^{(k)} = \sum_{l=1}^L w_f^{(k,l)} h_t^{(k,l)}$$

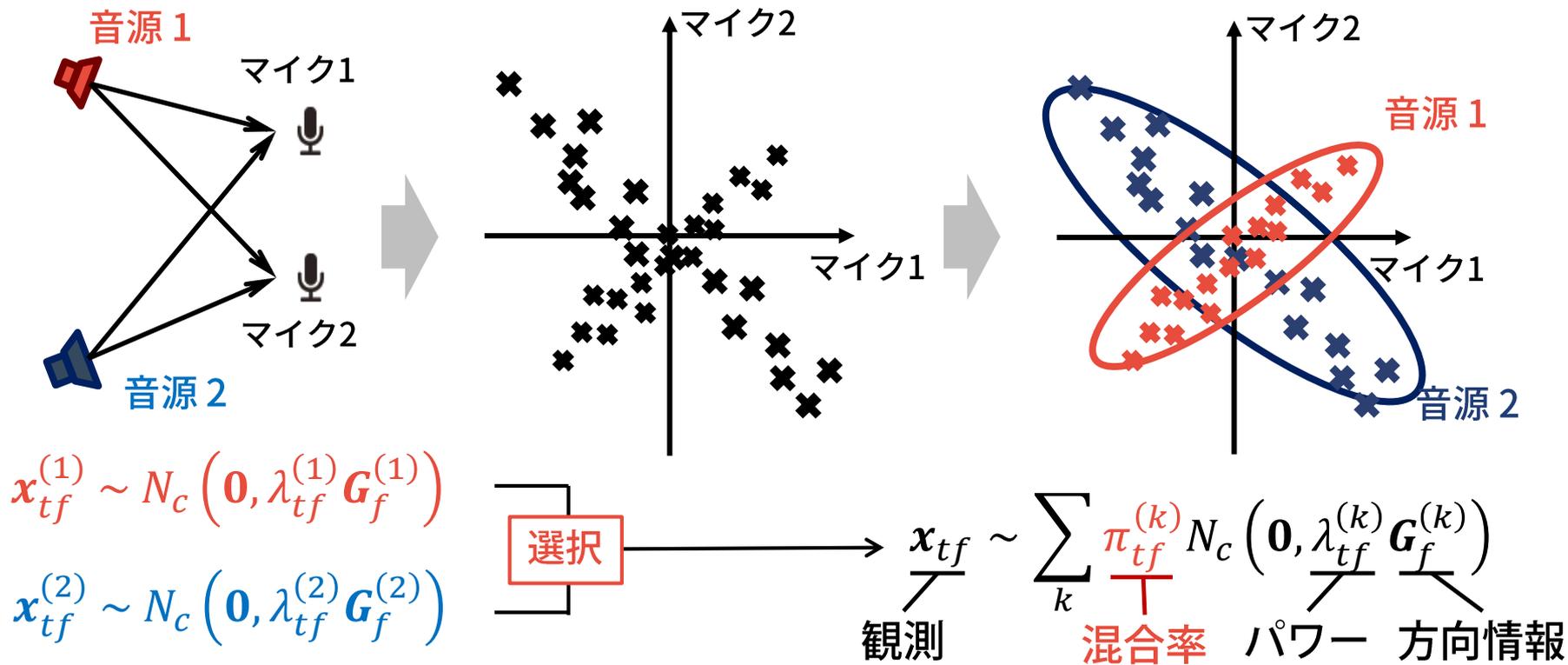
- 音源スペクトログラムを基底スペクトログラムの**パッチワーク**で表現
 - 音源 k の時間 t ・ 周波数 f におけるパワー λ_{tf}^k を、基底とアクティベーションの積のうちの**いずれか一つ**で表現



- 観測スペクトルを音源イメージの足し合わせで表現
 - 各時間 t ・周波数 f にはすべての音源が貢献



- 観測スペクトルを音源イメージの**パッチワーク**で表現
 - 各時間 t ・周波数 f には**いずれか一つ**のが貢献



四種類の統一的な確率モデル

- すべての組み合わせの定式化に成功 [Itakura 2016, 2017]
 - 周辺化ギブスサンプリングを用いたベイズ推論
 - すべての変数を交互にサンプリング

Σの位置に注意！

空間モデル

音源モデル

	因子モデル (factor)	混合モデル (mixture)
因子モデル (factor)	$N_c \left(\mathbf{0}, \sum_k \sum_l w_f^{(k,l)} h_t^{(k,l)} \mathbf{G}_f^{(k)} \right)$ 	$\sum_k \pi_{tf}^{(k)} N_c \left(\mathbf{0}, \sum_l w_f^{(k,l)} h_t^{(k,l)} \mathbf{G}_f^{(k)} \right)$ 
混合モデル (mixture)	$N_c \left(\mathbf{0}, \sum_k \sum_l \psi_{tf}^{(k,l)} w_f^{(k,l)} h_t^{(k,l)} \mathbf{G}_f^{(k)} \right)$ 	$\sum_k \sum_l \pi_{tf}^{(k)} \psi_{tf}^{(k,l)} N_c \left(\mathbf{0}, w_f^{(k,l)} h_t^{(k,l)} \mathbf{G}_f^{(k)} \right)$ 

- **マイクアレイ信号処理技術の基礎と最先端を紹介した**
 - **基礎的な事項**
 - ◆ 音の伝播のための線形システム
 - ◆ インパルス応答の意味・測定方法
 - **独立成分分析 (ICA)**
 - ◆ 主成分分析 (PCA) との違い
 - ◆ パーミュテーション問題・スケール問題
 - **時間周波数マスキング**
 - ◆ ノンパラメトリックベイズ拡張
 - **音源信号と観測信号の階層ベイズモデリング**
 - ◆ 音源モデル+空間モデル (それぞれ因子モデル or 混合モデル)

音楽情報処理や音響信号処理で一緒に研究してくださる方を探しています。
ぜひお問い合わせください。