# A Diagonal Plus Low-Rank Covariance Model for Computationally Efficient Source Separation

Antoine Liutkus (INRIA, France)
Kazuyoshi Yoshii (Kyoto University/RIKEN AIP, Japan)

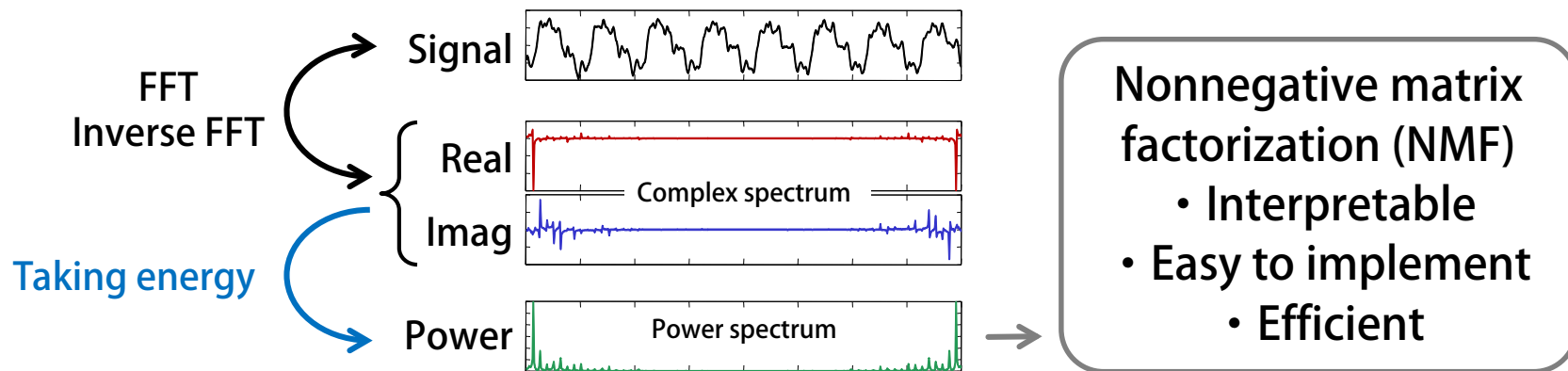# Outline

- We introduce positive semidefinite tensor factorization (PSDTF) based on the Log-Det divergence

  - A natural extension of nonnegative matrix factorization (NMF) based on the Itakura-Saito divergence

  - Estimation of locally-stationary Gaussian processes

- We propose a constrained version of LD-PSDTF for reducing computational complexity

  - Kernel matrices are restricted to diagonal + low-rank matrices

  - Woodbury formula is used for inversing kernel matrices

# Background

- ## Source separation is essential for various applications
  - Speech recognition and understanding
  - Automatic music transcription
- ## Phase information has not been used in most studies
  - The characteristics of sounds can be represented well in the power domain by discarding the phase information
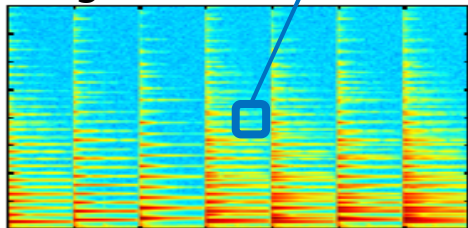  - The low-rankness and sparseness are useful clues

FFT
Inverse FFT

Signal

Real

Complex spectrum

Imag

Taking energy

Power

Power spectrum

Nonnegative matrix factorization (NMF)
- Interpretable
- Easy to implement
- Efficient

# Motivation

- ## Phase-aware source separation is promising
  - NMF can be extended based on additivity of complex spectra

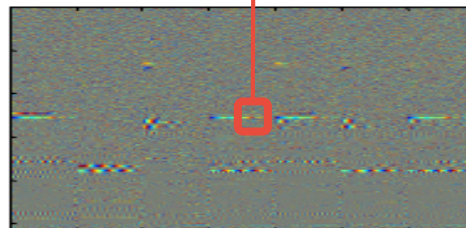|  | Frequency bins | Time frames |
|---|---|---|
| **Complex NMF** [Kameoka 2009] | Independent | Independent |
| **High Resolution NMF** [Badeau 2011] | Independent | Autoregressive |
| **PSDTF** [Yoshii 2013] | Correlated | Independent |

Additivity of time-domain signals

Complex value

$$x_{ft} = r_{ft}(\cos\theta_{ft} + i\sin\theta_{ft})$$

Magnitude $r$

Phase $\theta$
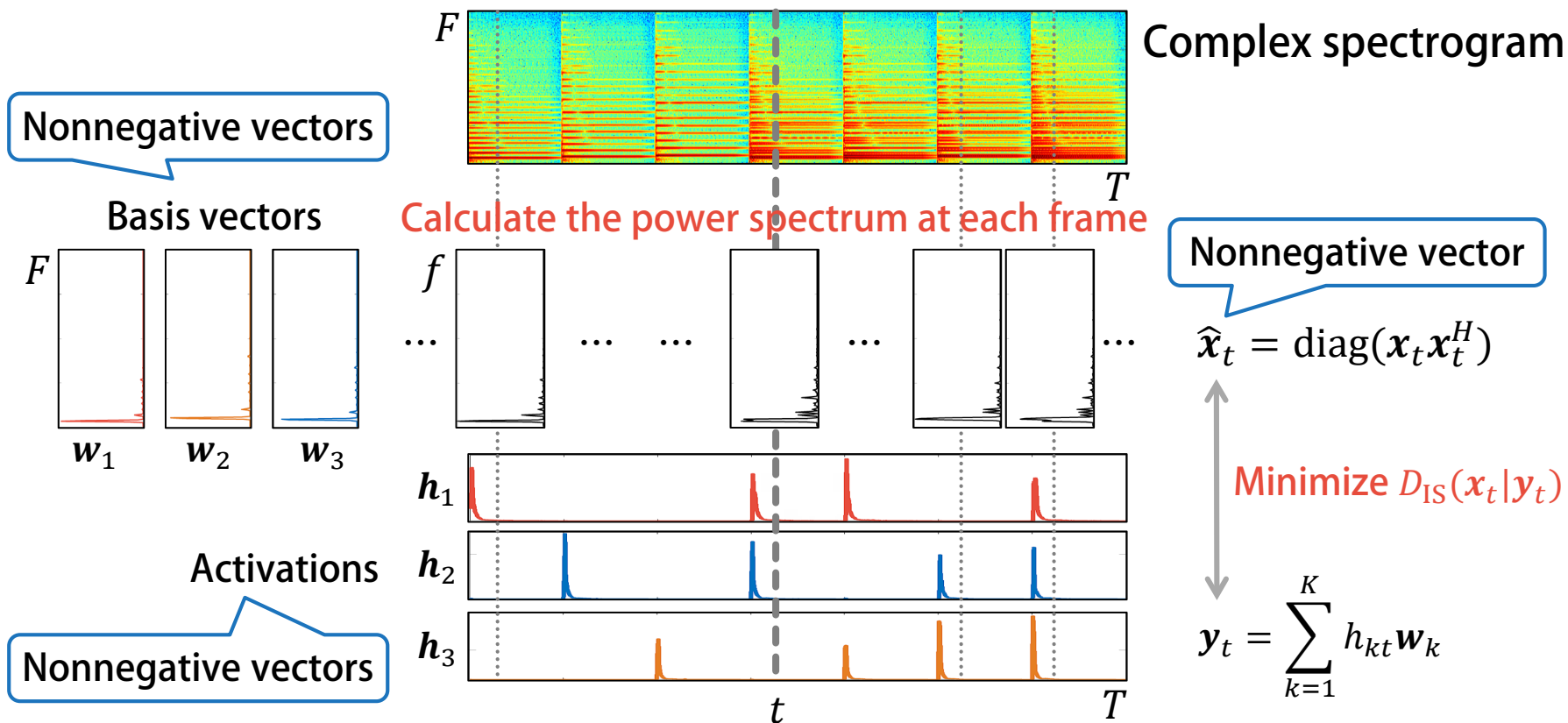


The values of magnitude and phase are not determined independently at frequency bins
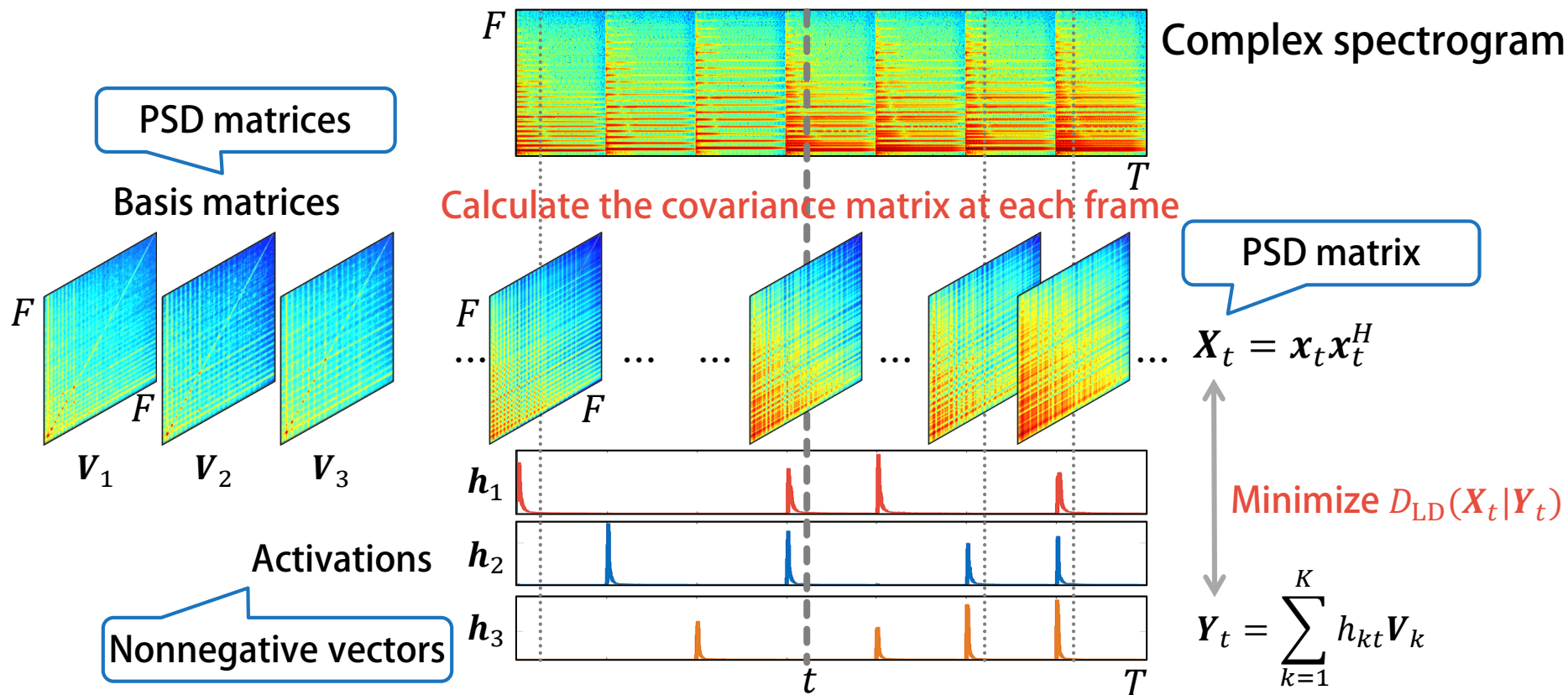
# Itakura-Saito NMF (IS-NMF) [Févotte 2009]

- Each observed <u>nonnegative vector</u> is approximated as the weighted sum of basis <u>nonnegative vectors</u>
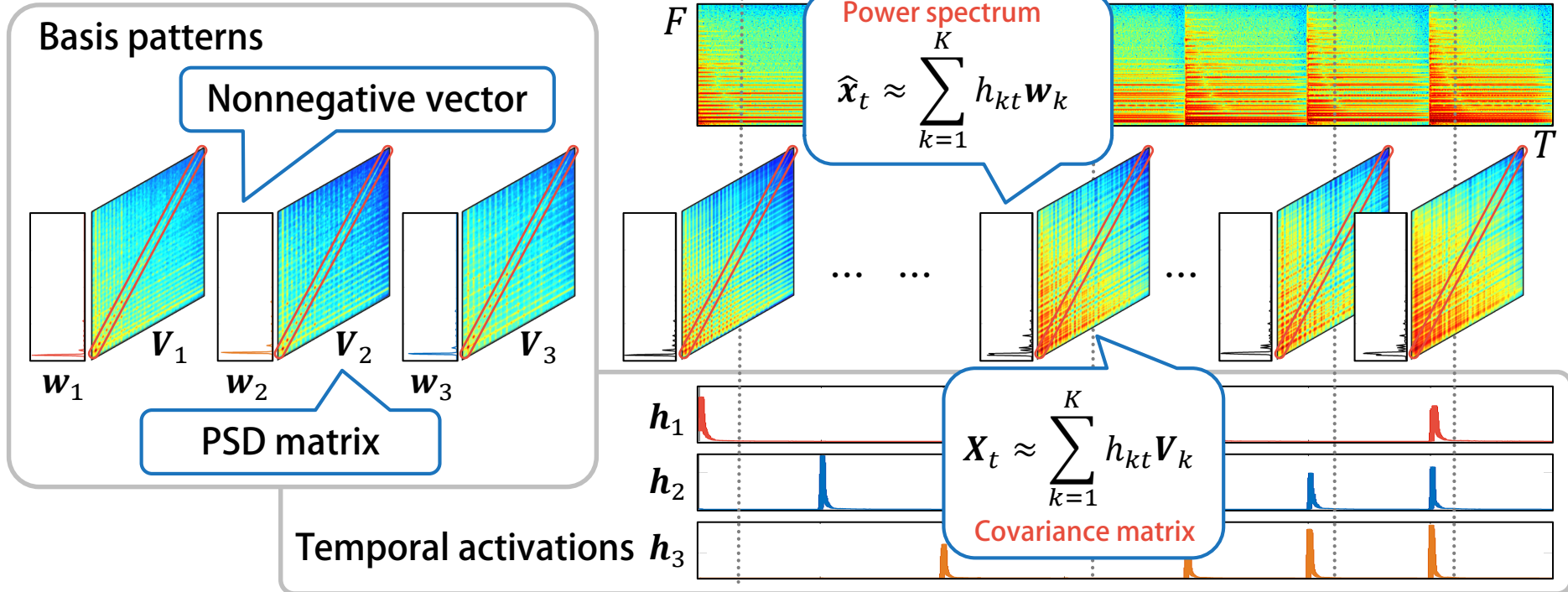
Complex spectrogram

Nonnegative vectors

Basis vectors

Calculate the power spectrum at each frame

Nonnegative vector

$$\widehat{x}_t = \mathrm{diag}(x_t x_t^H)$$

$w_1$   $w_2$   $w_3$

Minimize $D_{\mathrm{IS}}(x_t|y_t)$

Activations   $h_1$   $h_2$   $h_3$

Nonnegative vectors

$$y_t = \sum_{k=1}^{K} h_{kt} w_k$$

# Log-Det PSDTF (LD-PSDTF) [Yoshii 2013]

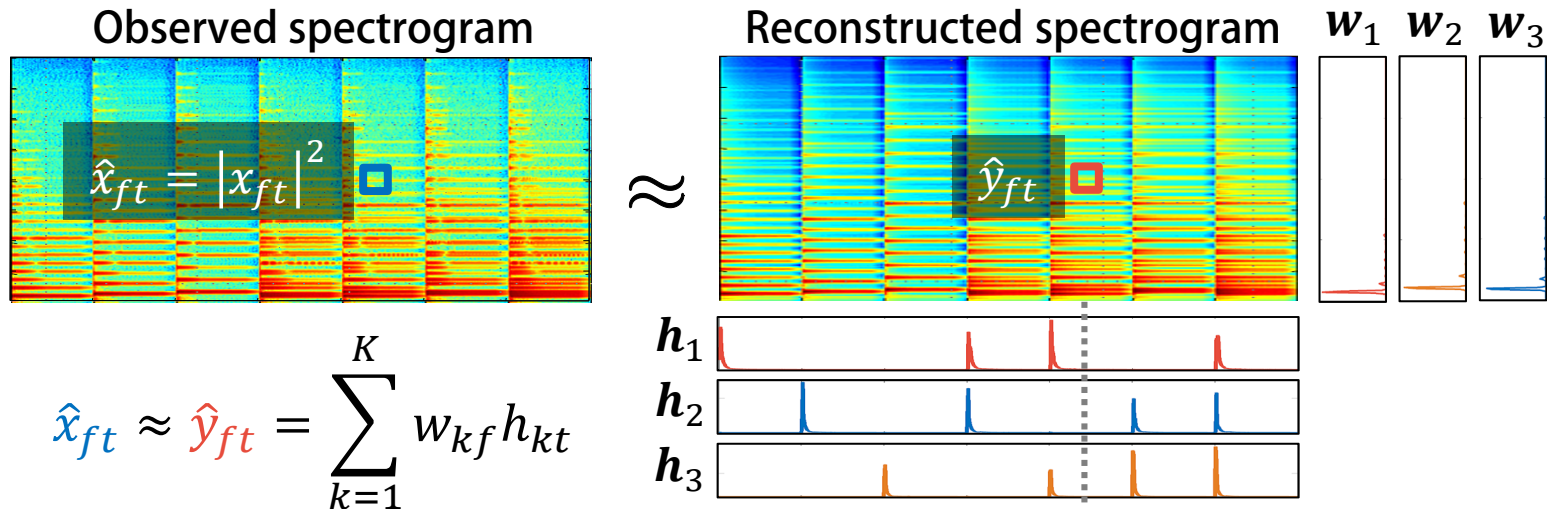- Each observed <u>pos. semidef. matrix</u> is approximated as the weighted sum of basis <u>pos. semidef. matrices</u>



Complex spectrogram

PSD matrices

Basis matrices

Calculate the covariance matrix at each frame

PSD matrix

$$X_t = x_t x_t^H$$

$V_1$  $V_2$  $V_3$

Activations  $h_1$  $h_2$  $h_3$

Nonnegative vectors

Minimize $D_{\mathrm{LD}}(X_t | Y_t)$

$$Y_t = \sum_{k=1}^{K} h_{kt} V_k$$

# IS-NMF vs LD-PSDTF

- ## LD-PSDTF is a natural extension of IS-NMF

  - Nonnegativity of scalars → Positive semidefinitenss of matrices
  - PSDTF reduces to NMF when all PSD matrices are diagonal



Basis patterns

Nonnegative vector

PSD matrix

$\boldsymbol{w}_1$ $\boldsymbol{V}_1$ $\boldsymbol{w}_2$ $\boldsymbol{V}_2$ $\boldsymbol{w}_3$ $\boldsymbol{V}_3$

$F$

Power spectrum

$$\widehat{\boldsymbol{x}}_t \approx \sum_{k=1}^{K} h_{kt} \boldsymbol{w}_k$$

$T$

… …  …

$\boldsymbol{h}_1$

$\boldsymbol{h}_2$

Temporal activations $\boldsymbol{h}_3$

$$\boldsymbol{X}_t \approx \sum_{k=1}^{K} h_{kt} \boldsymbol{V}_k$$

Covariance matrix

# Itakura-Saito NMF (IS-NMF) [Févotte 2009]

- ## NMF based on the Itakura-Saito divergence

  - The mixture spectrogram is approximated as a low-rank matrix
  - The number of sources $K$ should be specified in advance

Observed spectrogram

$$\hat{x}_{ft} = \left| x_{ft} \right|^2 \ \blacksquare$$

Reconstructed spectrogram

$$\hat{y}_{ft} \ \blacksquare$$

$w_1 \quad w_2 \quad w_3$

$\approx$

$$\hat{x}_{ft} \approx \hat{y}_{ft} = \sum_{k=1}^{K} w_{kf} h_{kt}$$

$h_1$

$h_2$

$h_3$

$$D_{\mathrm{IS}}\left( \hat{x}_{ft} \big| \hat{y}_{ft} \right) = -\log \frac{\hat{x}_{ft}}{\hat{y}_{ft}} + \frac{\hat{x}_{ft}}{\hat{y}_{ft}} - 1$$

Scale-invariant measure
$$D_{\mathrm{IS}}\left( \hat{x}_{ft} \big| \hat{y}_{ft} \right) = D_{\mathrm{IS}}\left( \alpha \hat{x}_{ft} \big| \alpha \hat{y}_{ft} \right)$$

# Log-Det PSDTF (LD-PSDTF) [Yoshii 2013]

- ## PSDTF based on the log-determinant divergence
  - The covariance matrix at each frame is approximated as the weighted sum of covariance matrices (basis matrices)

Observed spectrogram    Reconstructed spectrogram    $V_1$  $V_2$  $V_3$



$$X_t = x_t x_t^H$$

$$\approx$$

$$Y_t$$

$$X_t \approx Y_t = \sum_{k=1}^{K} V_k h_{kt}$$

$h_1$
$h_2$
$h_3$

$$D_{\mathrm{LD}}(X_t | Y_t) = -\log|X_t Y_t^{-1}| + \mathrm{tr}(X_t Y_t^{-1}) - F$$

Scale-invariant measure
$$D_{\mathrm{LD}}(X_t | Y_t) = D_{\mathrm{LD}}(\alpha X_t | \alpha Y_t)$$

# Probabilistic Formulation

- **The source signals are assumed to follow independent locally-stationary Gaussian processes in the time domain**

  - A mixture signal is the sum of multiple source signals

    | C | E | G | C+E | C+G | E+G | C+E+G |

    Assume the signals to be stationary in a short window

    $t$

    Mixture signal

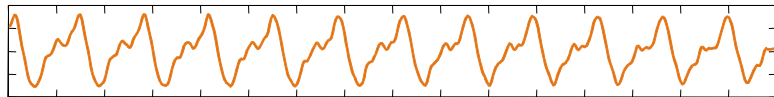    $$x_t = x_{1t} + x_{2t} + x_{3t}$$

    Source signal 1 (C)

    $$x_{1t} \sim N_c(0, h_{1t}V_1)$$

    Source signal 2 (E)

    $$x_{2t} \sim N_c(0, h_{2t}V_2)$$

    source signal 3 (G)

    $$x_{3t} \sim N_c(0, h_{3t}V_3)$$

# Mixing Process & Demixing Process

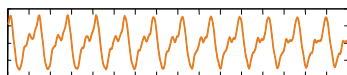- ## Sum of Gaussian variables → Gaussian variable



$x_{1t} \sim N_c(0, Y_{1t})$
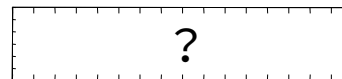
$x_{2t} \sim N_c(0, Y_{2t})$

$x_{3t} \sim N_c(0, Y_{3t})$

?

$x_t = x_{1t} + x_{2t} + x_{3t}$

$\sim N_c(0, Y_{1t} + Y_{3t} + Y_{3t} = Y_t)$

- ## Gaussian variable → Sum of Gaussian variables

?

$x_{1t} \sim N_c(0, Y_{1t})$

?

$x_{2t} \sim N_c(0, Y_{2t})$

?

$x_{3t} \sim N_c(0, Y_{3t})$



$x_t = x_{1t} + x_{2t} + x_{3t}$

$\sim N_c(0, Y_{1t} + Y_{3t} + Y_{3t} = Y_t)$

$x_{1t}|x_t \sim N_c(Y_{1t}Y_t^{-1}x_t, Y_{1t} - Y_{1t}Y_t^{-1}Y_{1t})$

$x_{2t}|x_t \sim N_c(Y_{2t}Y_t^{-1}x_t, Y_{2t} - Y_{2t}Y_t^{-1}Y_{2t})$

$x_{3t}|x_t \sim N_c(Y_{3t}Y_t^{-1}x_t, Y_{3t} - Y_{3t}Y_t^{-1}Y_{3t})$

All the frequency components of each source spectrum can be estimated jointly via Wiener filtering

# Maximum Likelihood Estimation

- We aim to estimate $H, V$ that maximizes the likelihood

Observed complex spectrogram



$x_t$

$$x_t \sim N_c\left(\mathbf{0}, \sum_{k=1}^{K} h_{kt} V_k\right) \longrightarrow \text{Maximize}$$

**Observed covariance matrix**

$$X_t = x_t x_t^H$$

$\longleftrightarrow$

**Approx. covariance matrix**

$$Y_t = \sum_{k=1}^{K} h_{kt} V_k$$

**Gaussian log-likelihood**

$$\log p(X_t | Y_t) = -\frac{1}{2}\log|Y_t| - \frac{1}{2}\mathrm{tr}(X_t Y_t^{-1}) \longrightarrow \text{Maximize}$$

**Log-Det divergence**

$$D(X_t | Y_t) = -\log|X_t Y_t^{-1}| + \mathrm{tr}(X_t Y_t^{-1}) - F \longrightarrow \text{Minimize}$$

Equivalent!

# Generalized EM Algorithm (Proposed)

- Iteratively update latent sources and parameters
  - Expectation step
    - Calculate covariance matrices $\boldsymbol{Y}_{kt} = h_{kt}\boldsymbol{V}_k \quad \boldsymbol{Y}_t = \sum_{k=1}^{K} \boldsymbol{Y}_{kt}$
    - Calculate posteriors of source spectra
    $$\boldsymbol{x}_{kt}|\boldsymbol{x}_t \sim N_c\big(\underbrace{\boldsymbol{Y}_{kt}\boldsymbol{Y}_t^{-1}\boldsymbol{x}_t}_{\mathrm{E}[\boldsymbol{x}_{kt}]}, \underbrace{\boldsymbol{Y}_{kt} - \boldsymbol{Y}_{kt}\boldsymbol{Y}_t^{-1}\boldsymbol{Y}_{kt}}_{\mathrm{V}[\boldsymbol{x}_{kt}]}\big)$$
    - Calculate second-order statistics
    $$\mathrm{E}\big[\boldsymbol{x}_{kt}\boldsymbol{x}_{kt}^H\big] = \mathrm{E}[\boldsymbol{x}_{kt}]\mathrm{E}[\boldsymbol{x}_{kt}^H] + \mathrm{V}[\boldsymbol{x}_{kt}]$$

    > IS-NMF: $O(KTF)$
    > LD-PSDTF: $O(KTF^3)$

  - Maximization step
    - Update parameters (depend on each other)
    $$h_{kt} \leftarrow \frac{\mathrm{tr}\big(\boldsymbol{V}_k^{-1}\mathrm{E}\big[\boldsymbol{x}_{kt}\boldsymbol{x}_{kt}^H\big]\big)}{F} \qquad \boldsymbol{V}_k \leftarrow \frac{\sum_{t=1}^{T} h_{kt}^{-1}\mathrm{E}\big[\boldsymbol{x}_{kt}\boldsymbol{x}_{kt}^H\big]}{T}$$

# Computational Bottleneck

- Inversion of big matrices is computationally prohibitive

  - E step: updating source spectra

  $$\mathrm{E}\big[\boldsymbol{x}_{kt}\boldsymbol{x}_{kt}^{H}\big] = \boldsymbol{Y}_{kt}\textcolor{red}{\boldsymbol{Y}_{t}^{-1}}\boldsymbol{x}_{t} + \boldsymbol{Y}_{kt} - \boldsymbol{Y}_{kt}\textcolor{red}{\boldsymbol{Y}_{t}^{-1}}\boldsymbol{Y}_{kt}$$

  - M step: updating parameters

  $$h_{kt} \leftarrow \frac{\mathrm{tr}\big(\textcolor{red}{\boldsymbol{V}_{k}^{-1}}\mathrm{E}\big[\boldsymbol{x}_{kt}\boldsymbol{x}_{kt}^{H}\big]\big)}{F} \qquad \boldsymbol{V}_{k} \leftarrow \frac{\sum_{t=1}^{T} h_{kt}^{-1}\mathrm{E}\big[\boldsymbol{x}_{kt}\boldsymbol{x}_{kt}^{H}\big]}{T}$$
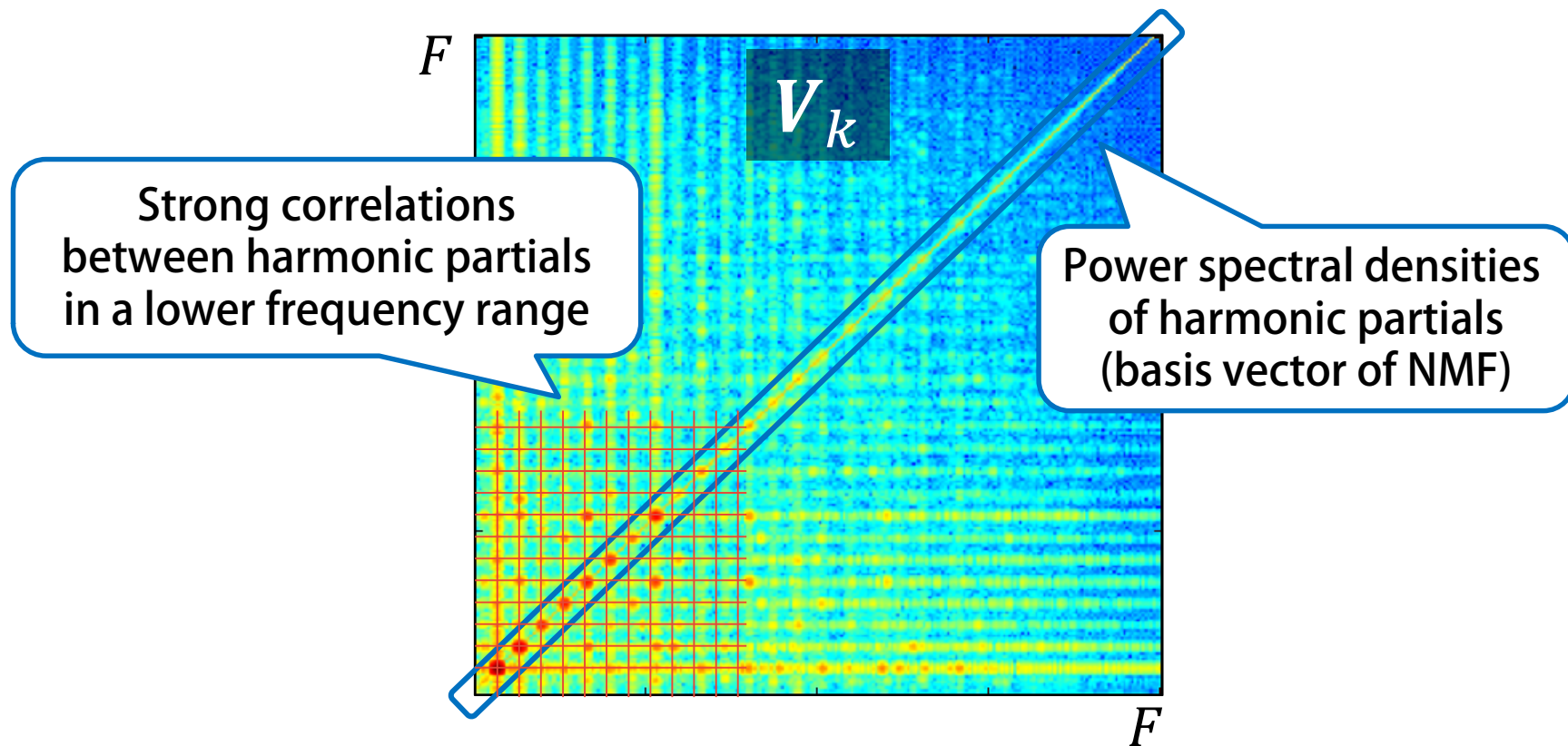
The inverse matrices $\boldsymbol{Y}_{t}^{-1}$ and $\boldsymbol{V}_{k}^{-1} \in \mathrm{C}^{F \times F}$ are required: $O(F^3)$

How to calculate these inversions
in a more efficient manner?

# Covariance Matrix Revisited

- Basis covariance matrices have diagonal + grid patterns
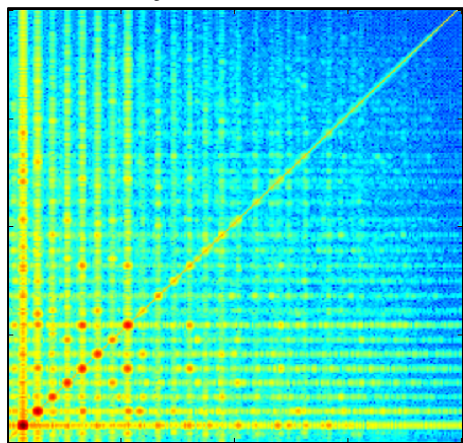  - Especially for complex spectra with harmonic structures



$F$

$V_k$

Strong correlations between harmonic partials in a lower frequency range

Power spectral densities of harmonic partials (basis vector of NMF)

$F$

# Covariance Approximation (Proposed)

- Each $V_k$ is approximated as a diagonal + low-rank matrix
  - The rank $N$ can be around the number of harmonic partials
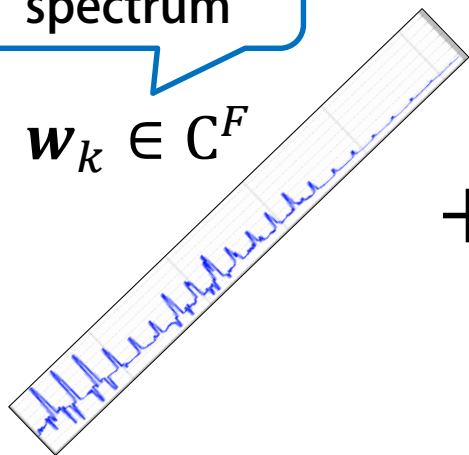
$$V_k = [w_k] + L_k[s_k]L_k^H$$

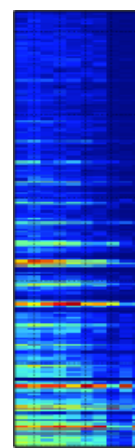Basis covariance matrix

$$V_k \in \mathrm{C}^{F \times F}$$
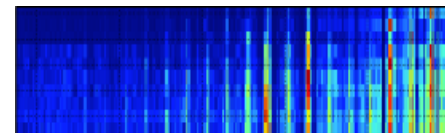
Basis power spectrum

$$w_k \in \mathrm{C}^F$$

$$s_k \in \mathrm{C}^N$$

$$L_k^H \in \mathrm{C}^{N \times F}$$

$$= \qquad + \qquad L_k \in \mathrm{C}^{F \times N}$$

# EM Algorithm Revisited

- The inversion of big matrices are required

  - E step: updating source spectra

  $$\mathrm{E}\left[\boldsymbol{x}_{kt}\boldsymbol{x}_{kt}^{H}\right] = \boldsymbol{Y}_{kt}\boldsymbol{Y}_{t}^{-1}\boldsymbol{x}_{t} + \boldsymbol{Y}_{kt} - \boldsymbol{Y}_{kt}\boldsymbol{Y}_{t}^{-1}\boldsymbol{Y}_{kt}$$

  - M step: updating parameters

  $$h_{kt} \leftarrow \frac{\mathrm{tr}\left(\boldsymbol{V}_{k}^{-1}\mathrm{E}\left[\boldsymbol{x}_{kt}\boldsymbol{x}_{kt}^{H}\right]\right)}{F} \qquad \boldsymbol{V}_{k} \leftarrow \frac{\sum_{t=1}^{T}h_{kt}^{-1}\mathrm{E}\left[\boldsymbol{x}_{kt}\boldsymbol{x}_{kt}^{H}\right]}{T}$$

  $$\begin{cases} \boldsymbol{V}_{k} = [\boldsymbol{w}_{k}] + \boldsymbol{L}_{k}[\boldsymbol{s}_{k}]\boldsymbol{L}_{k}^{H} \\ \boldsymbol{Y}_{t} = \sum_{k=1}^{K}h_{kt}\boldsymbol{V}_{k} = \sum_{k=1}^{K}h_{kt}[\boldsymbol{w}_{k}] + \sum_{k=1}^{K}h_{kt}\boldsymbol{L}_{k}[\boldsymbol{s}_{k}]\boldsymbol{L}_{k}^{H} \end{cases}$$
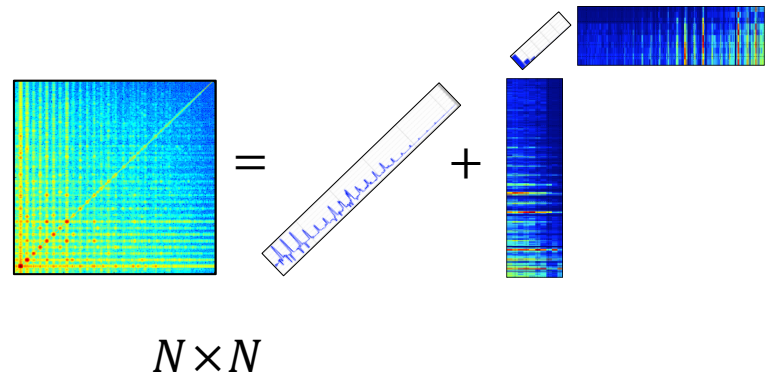
  Each term can be inverted efficiently

# Efficient Matrix Inversion

- Use Woodbury formula for covariance matrices

$$(A + UCV)^{-1} = A^{-1} - A^{-1}U\boxed{(C^{-1} + VA^{-1}U)^{-1}}VA^{-1}$$

This formula is useful when $A$ and $C$ can be inverted efficiently

Diagonal matrices!

$$\underset{F \times F}{V_k} = \underset{F \times F}{[w_k]} + \underset{\substack{F \times N \ N \times N \ N \times F \\ (N \ll F)}}{L_k [s_k] L_k^H}$$



$N \times N$

$$\underset{F \times F}{V_k^{-1}} = \underset{F \times F}{[w_k]^{-1}} - \underset{\substack{F \times F \quad F \times N}}{[w_k]^{-1}L_k} \boxed{\underset{\substack{N \times N \quad N \times F \ F \times F \quad F \times N}}{\left([s_k]^{-1} + L_k^H[w_k]^{-1}L_k\right)^{-1}}} \underset{\substack{N \times F \ F \times F}}{L_k^H[w_k]^{-1}}$$

Inversion of a compact matrix!

# Recursive Matrix Inversion

- Use Woodbury formula in a recursive manner

$$Y_t = \sum_{k=1}^{K} h_{kt}[\boldsymbol{w}_k] + \sum_{k=1}^{\textcircled{K}} h_{kt}\boldsymbol{L}_k[\textcolor{red}{\boldsymbol{s}_k}]\boldsymbol{L}_k^H \implies \boldsymbol{Y}_t^{-1}$$

$$\boldsymbol{Y}_t^{(p)} \stackrel{\text{def}}{=} \sum_{k=1}^{K} h_{kt}[\boldsymbol{w}_k] + \sum_{k=1}^{\textcircled{p}} h_{kt}\boldsymbol{L}_k[\textcolor{red}{\boldsymbol{s}_k}]\boldsymbol{L}_k^H = \boldsymbol{Y}_t^{(p-1)} + h_{pt}\boldsymbol{L}_p[\textcolor{red}{\boldsymbol{s}_p}]\boldsymbol{L}_p^H$$

$$\left(\boldsymbol{Y}_t^{(p)}\right)^{-1} = \left(\boldsymbol{Y}_t^{(p-1)}\right)^{-1}$$

$$- \left(\boldsymbol{Y}_t^{(p-1)}\right)^{-1} \boldsymbol{L}_p \overset{N\times N}{\boxed{\left(h_{pt}^{-1}[\textcolor{red}{\boldsymbol{s}_p}]^{-1} + \boldsymbol{L}_p^H \left(\boldsymbol{Y}_t^{(p-1)}\right)^{-1} \boldsymbol{L}_p\right)^{-1}}} \boldsymbol{L}_p^H \left(\boldsymbol{Y}_t^{(p-1)}\right)^{-1}$$

Recurrence formula starting at $\boldsymbol{Y}_t^{(0)} = \sum_{k=1}^{K} h_{kt}[\boldsymbol{w}_k]$ (NMF)
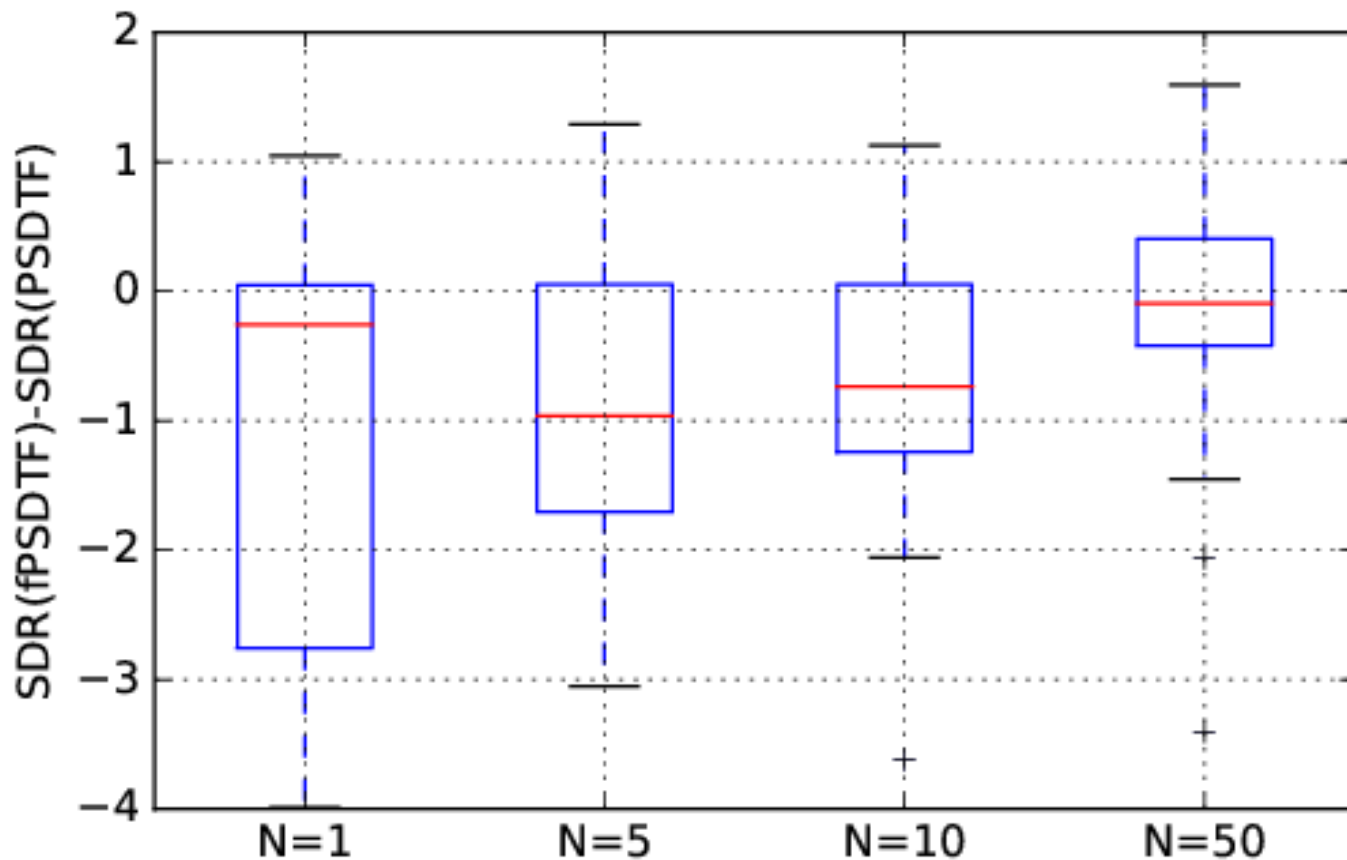
# Evaluation

- **Separation performance vs covariance approximation**
  - Synthesize a mixture signal sampled at 16 [kHz]
    - $K = 3$ (C4, E4, G4, piano) $\cdot$ $F = 256, T = 840$
  - Test "fast" PSDTF wih $N = 0$ (NMF), $1, 5, 10, 50, 256$ (PSDTF)
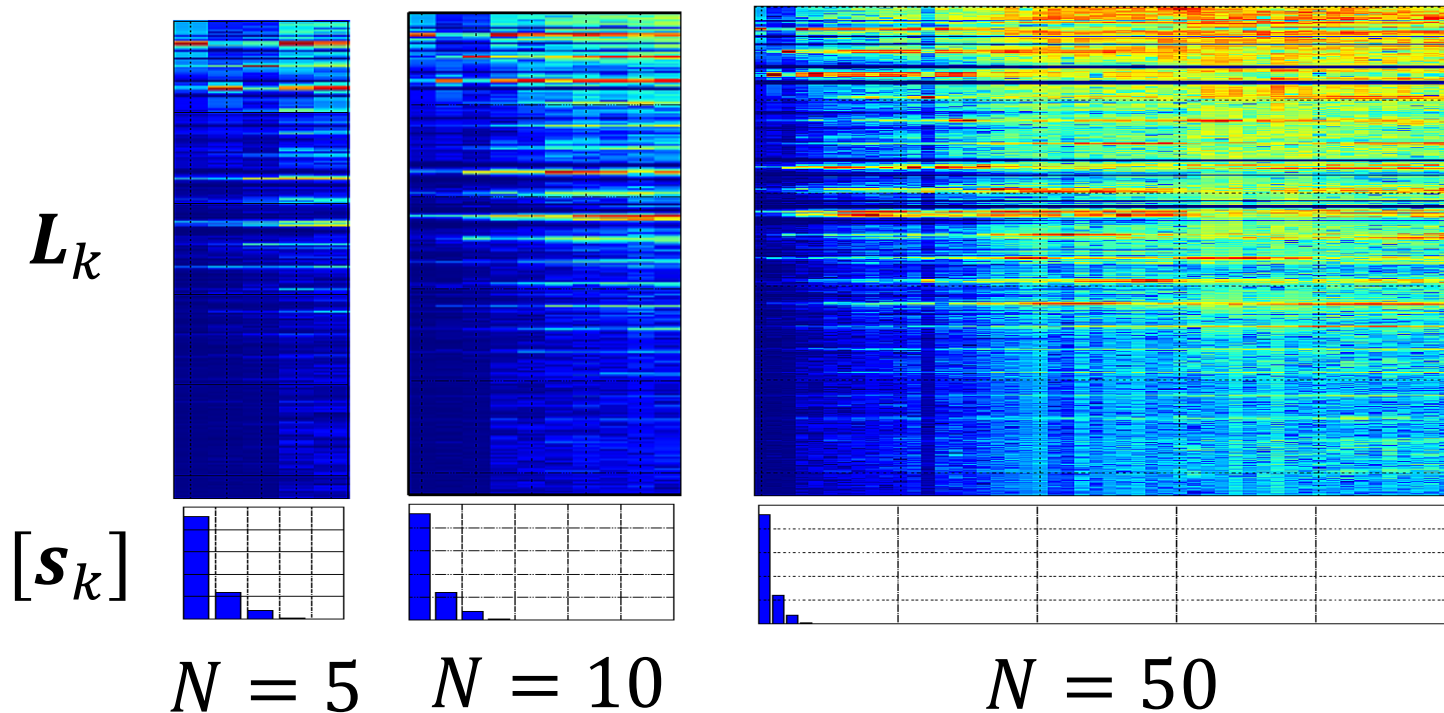  - Use BSS Eval Toolbox [Vincent2006]

# Separation Performance

- PSDTF ($N = 50$) was comparable with PSDTF ($N = 256$)

# Estimated Results

- The off-diagonal elements of each $V_k$ (inter-frequency correlations) can be approximated by a low-rank matrix
  - A limited number of eigenvalues are significantly larger than 0



$L_k$

$[s_k]$

$N = 5$   $N = 10$   $N = 50$

# Conclusion

- We introduced positive semidefinite tensor factorization (PSDTF) based on the Log-Det divergence

  - A natural extension of nonnegative matrix factorization (NMF) based on the Itakura-Saito divergence

  - Estimation of locally-stationary Gaussian processes

- We proposed a constrained version of LD-PSDTF for reducing computational complexity

  - Kernel matrices are restricted to diagonal + low-rank matrices

  - Woodbury formula is used for inversing kernel matrices (in a recursive manner)

$$\boldsymbol{V}_k = [\boldsymbol{w}_k] + \boldsymbol{L}_k [\boldsymbol{s}_k] \boldsymbol{L}_k^H$$