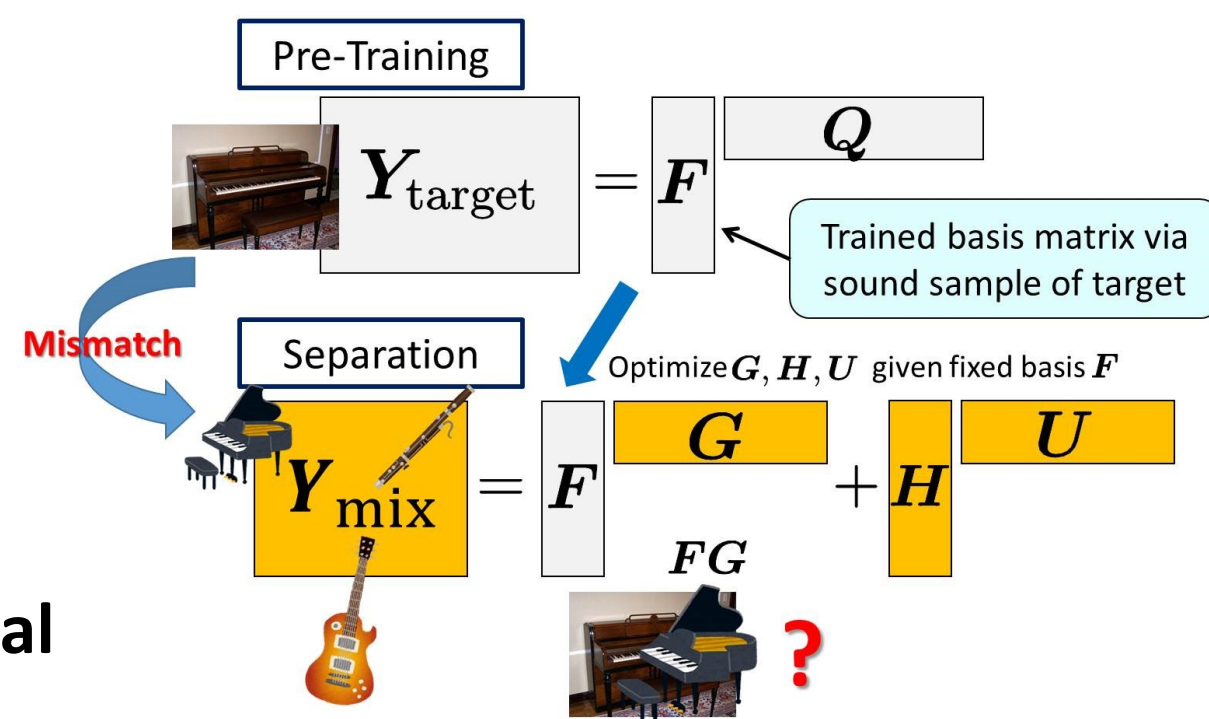


Music Signal Separation Using Supervised NMF with All-Pole-Model-Based Discriminative Basis Deformation

H. Nakajima (The University of Tokyo: UTokyo), D. Kitamura (SOKENDAI), N. Takamune (UTokyo), S.Koyama (UTokyo), H. Saruwatari (UTokyo), N. Ono (NII) , Y. Takahashi (Yamaha Corp.), K. Kondo (Yamaha Corp.)

1. INTRODUCTION

- Nonnegative matrix factorization (NMF)* can be used to decompose a spectrogram into BASIS matrix and ACTIVATION matrix.
- Supervised NMF (SNMF)* can extract the target source by using pre-trained basis matrix F in advance.
- DRAWBACK of SNMF**
 - Timbre mismatch between pre-trained basis and actual target signal causes big error in source separation.



RESEARCH PURPOSE

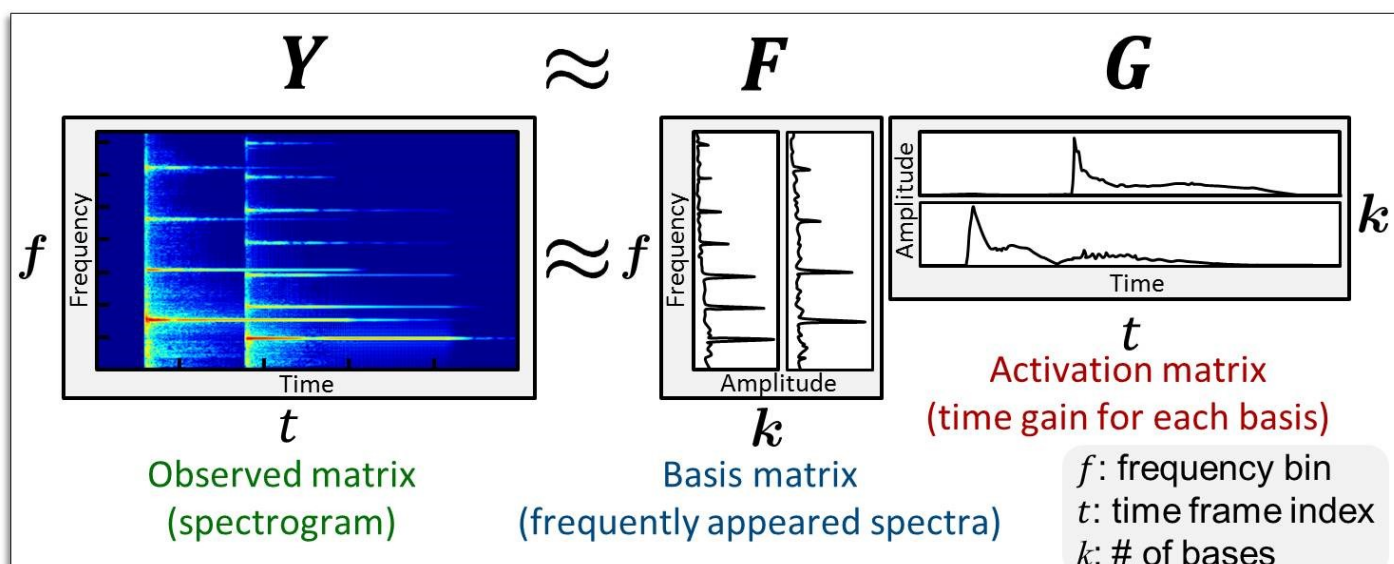
- Introduction of **All-Pole-Model-Based Basis Deformation**
- Discriminative Basis Training** for obtaining appropriate degree of freedom of deformation

2. CONVENTIONAL METHODS

NMF and SNMF [1]

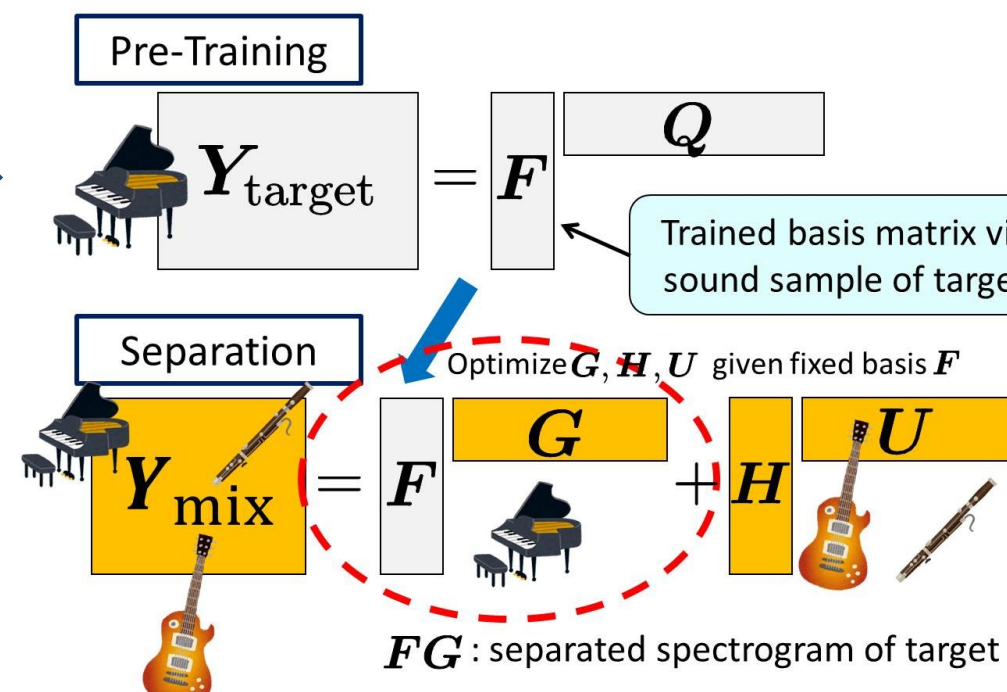
NMF:

Sparse representation using nonnegative matrices F and G



SNMF:

Using pre-trained basis matrix F , we can decompose $Y \cong FG + HU$ to get FG (target source spectrogram).



SNMF with Additive Basis Deformation (SNMF-ABD) [2]

The method deforms the original basis F to be $F + D$ to adapt the actual target spectrogram in Y_{mix} .

Deformation model

$$Y_{\text{mix}} \approx (F + D)G + HU$$

- ✗ Difficult to optimize because D is NOT limited within nonnegative.
- ✗ Hard to simultaneously perform separation and deformation.

3. PROPOSED METHOD FOR BASIS DEFORMATION

Statistical Postfilter: Generalized MMSE Short-Time Spectral Amplitude (MMSE-STSA) Estimator [3,4]

High accuracy NMF postfilter based on non-target spectral variance and a priori statistical target model.

$$Y = J \circ Y_{\text{mix}}$$

\circ : Hadamard product

Assuming target amplitude spectrum Y_{target} obeys chi distribution, we solve $\min_j E[(Y_{\text{target}} - Y)^2]$

$$J_{\omega,t} = \frac{\sqrt{v_{\omega,t}}}{\tilde{\gamma}_{\omega,t}} \cdot \left(\frac{\Gamma(\rho + 0.5)}{\Gamma(\rho)} \cdot \frac{\Phi(0.5 - \rho, 1, -v_{\omega,t})}{\Phi(1 - \rho, 1, -v_{\omega,t})} \right)^{1/\beta}$$

$\Phi(a, b; k)$: confluent hyper-geometric function

$\Gamma()$: gamma function

ϵ, γ : a priori and a posteriori SNR and $v_{\omega,t} = \tilde{\gamma}_{\omega,t} \tilde{\epsilon}_{\omega,t} (1 + \tilde{\epsilon}_{\omega,t})^{-1}$

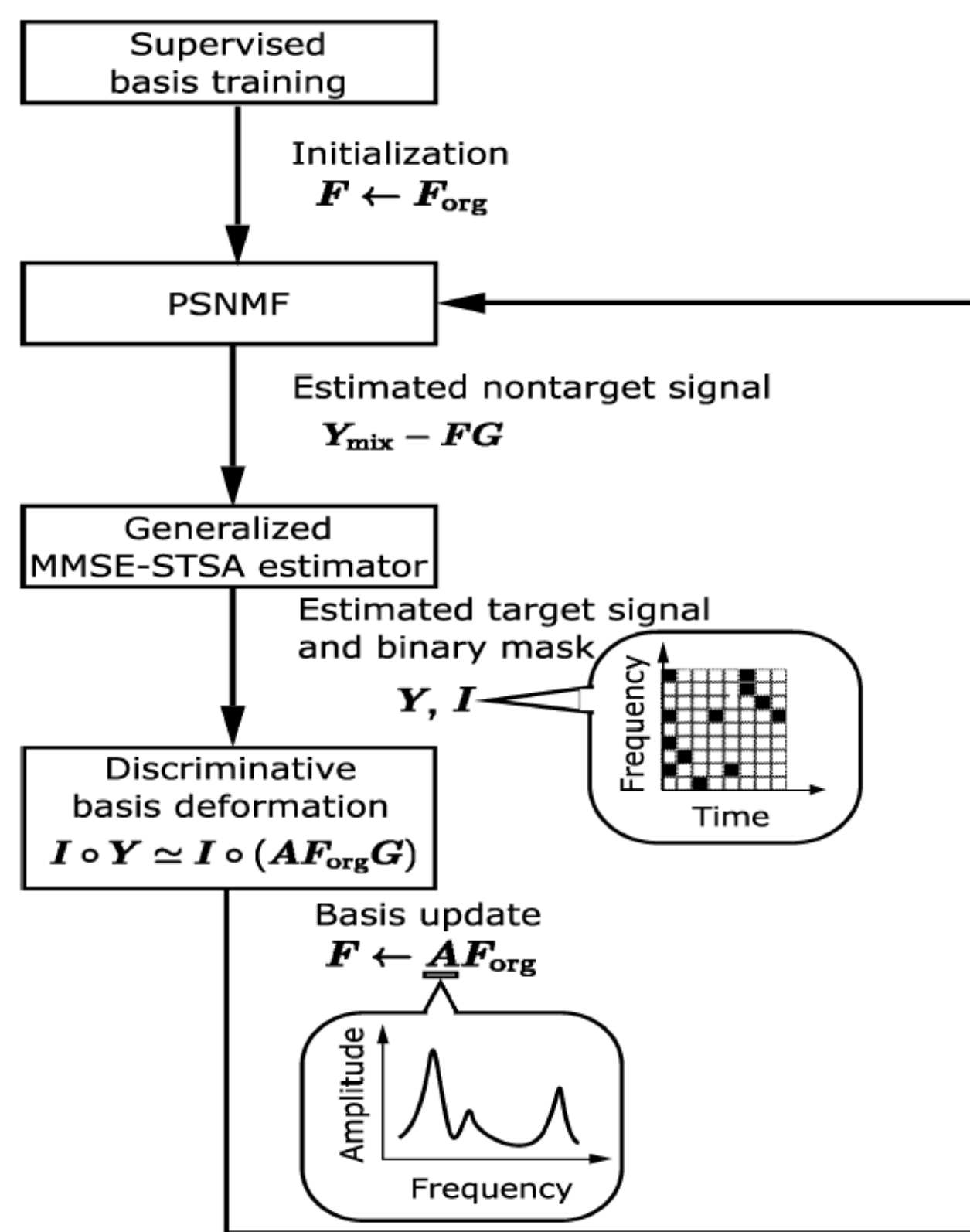
ρ : shape parameter of chi distribution

Thanks to the a priori target statistical model, we can obtain higher separation result rather than SNMF only.

J can be used as the measure of seldom overlapping components among target and non-target signals.
⇒ we define a matrix I as the binary version of J for sampling convincing components

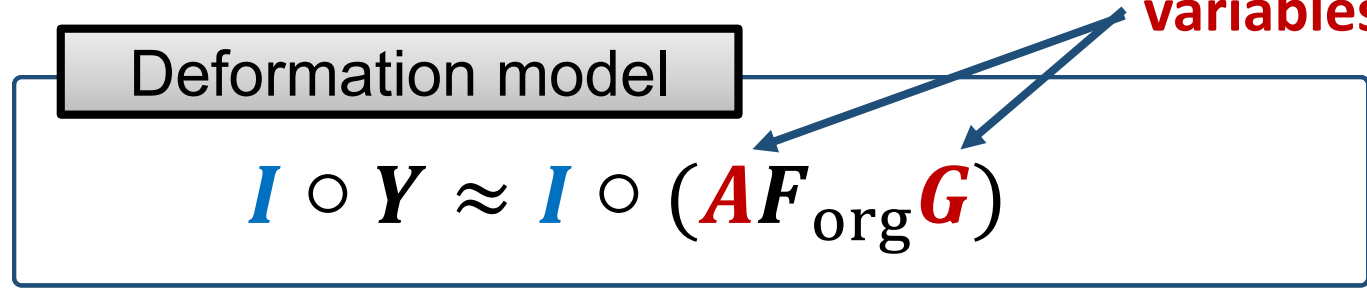
γ can be calculated by NMF output $Y_{\text{mix}} - FG$

Basic Idea and System Overview



Motivation:

- Simultaneous processing for separation and basis deformation results in bad side effect each other.
- We want to separate the whole process into 2 parts: **deformation** and **separation**.



I : Binary masking to sample convincing target components

Y : Output of Generalized MMSE-STSA estimator

A : Diagonal matrix with All-Pole-Model spectrum gain

F_{org} : Original basis matrix

G : Activation matrix

Smoothed spectral envelope to reduce mismatch between the sampled target spectra and F_{org}

Cost Function and Parameter Update Rule

Cost function to optimize all-pole-model A using Y and I (KL-divergence-based).

$$\mathcal{J} = \sum_{\omega,t} i_{\omega,t} \left\{ -y_{\omega,t} + \frac{\sum_k f_{\omega,k} g_{k,t}}{|A_{\omega}|} + y_{\omega,t} \log \frac{y_{\omega,t}}{\sum_k f_{\omega,k} g_{k,t} / |A_{\omega}|} \right\}$$

where $A_{\omega} = 1 - \sum_{k=1}^p \alpha_k \exp(-\pi j k \frac{\omega}{\Omega})$

$i_{\omega,t}$: element of I $y_{\omega,t}$: element of Y

$f_{\omega,k}$: element of F_{org} $g_{k,t}$: element of G

Ω : Nyquist frequency

Update Rule

$$\alpha \leftarrow R^{-1}r \quad \alpha \text{ is a vector with entry } \alpha_k$$

$$R_{k,q} = \sum_{\omega,t} \left[i_{\omega,t} \left(\sum_k f_{\omega,k} g_{k,t} \frac{1}{|A_{\omega}|^3} + y_{\omega,t} \frac{1}{2|A_{\omega}|^2} \right) \left(\exp(-2\pi j \frac{\omega}{\Omega} (k-q)) + \exp(2\pi j \frac{\omega}{\Omega} (k-q)) \right) \right]$$

$$r_q = \sum_{\omega,t} i_{\omega,t} \left[\left(\sum_k f_{\omega,k} g_{k,t} \frac{1}{|A_{\omega}|^3} + y_{\omega,t} \frac{1}{2|A_{\omega}|^2} \right) \left(\exp(-2\pi j \frac{\omega}{\Omega} q) + \exp(2\pi j \frac{\omega}{\Omega} q) \right) - \frac{3}{|A_{\omega}|^2} \sum_k f_{\omega,k} g_{k,t} \text{Re} \left[\frac{A_{\omega}^*}{|A_{\omega}|} \exp(-2\pi j \frac{\omega}{\Omega} q) \right] \right]$$

$$g_{k,t} \leftarrow g_{k,t} \frac{\sum_{\omega} i_{\omega,t} y_{\omega,t} f_{\omega,k} / (\sum_{\omega} f_{\omega,k} g_{k,t})}{\sum_{\omega} i_{\omega,t} f_{\omega,k} / |A_{\omega}|}$$

4. DISCRIMINATIVE BASIS TRAINING

Problem in determining appropriate order of all-pole model

- ✗ Small order leads to insufficient deformation (basis mismatch cannot be resolved).
- ✗ Large order results in exceeding deformation that wrongly represents other (non-target) signal.

We should solve **Bilevel Optimization** that can supply the appropriate solution holding both conditions: (a) spectral mismatch becomes as small as possible, and (b) mixture can be modeled as accurate as possible using HU .

$$A = \arg \min_A \mathcal{D}_{\text{KL}} \left(I \circ Y | I \circ (A F G_1) \right) \quad \text{s.t. } G_1$$
$$= \arg \min_{G,H,U} \mathcal{D}_{P_KL} \left(I \circ Y_{\text{mix}} | I \circ (A F G + H U) \right)$$

Since the bilevel optimization is generally hard to solve, we introduce an **approximated iterative parameter update rule** as follows.

Step 1 : Initialization

$$A_s = \arg \min_{A, G} \mathcal{D}_{\text{KL}} \left(I \circ Y | I \circ (A F G) \right).$$

Step 2 : Modeling of Mixture Y_{mix}

$$G_s = \arg \min_{G, H, U} \mathcal{D}_{P_KL} \left(I \circ Y_{\text{mix}} | I \circ (A_s F G + H U) \right).$$

Step 3 : Modeling of Target Y

$$A_s = \arg \min_A \mathcal{D}_{\text{KL}} \left(I \circ Y | I \circ (A F G_s) \right).$$

Return to Step 2

5. EXPERIMENTAL EVALUATION

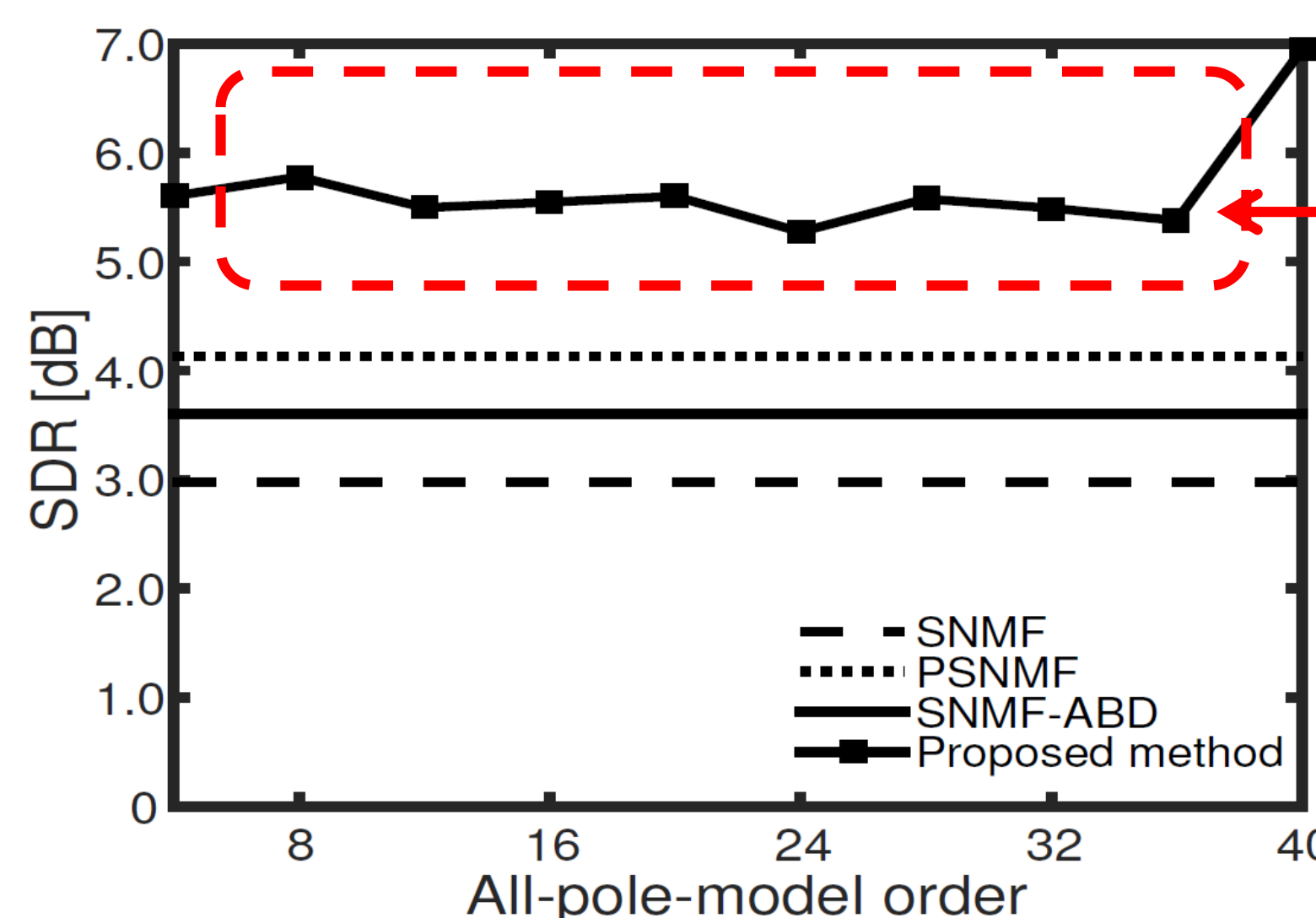
Experimental Condition

Type of instruments	Oboe (Ob.), Piano (Pf.), Trombone (Tb.)
Training sample (MIDI)	2 octave notes generated by Garritan Professional Orchestra
Target sound (MIDI)	Generated by Microsoft GS Wavetable SW Synth (*)
Sampling frequency	44100 Hz
FFT length	4096 points (about 92 ms)
Shift length	512 points (about 12 ms)
Number of bases	Target: 100, Non-target: 30
Iteration	Training : 1000, Separation: 1000
Evaluation score	Signal-to-Distortion Ratio (SDR)
Comparison	SNMF, Penalized SNMF (PSNMF) [5], SNMF-ABD [2]

*Music score used in this experiment



Results



Proposed method outperforms the conventional NMF methods in almost all the cases.

Example of SDR for separating Pf. from the mixture of Pf. and Ob.

MAXIMUM VALUE OF SDR IN EACH MIXTURE [dB]				
	SNMF	PSNMF	SNMF-ABD	Proposed method
Ob. & Pf.	7.6	6.7	8.1	7.1
Ob. & Tb.	1.5	2.4	2.6	3.0
Pf. & Ob.	3.0	4.1	3.6	7.0
Pf. & Tb.	1.9	3.1	3.2	5.0
Tb. & Ob.	-0.6	0.7	0.2	2.6
Tb. & Pf.	1.8	2.9	2.6	4.5
Average	2.5	3.3	3.4	4.9

REFERENCES

- [1] D. D. Lee, et al., *Proc. Advances in Neural information Processing Systems*, 2001.
- [2] D. Kitamura, et al., *Proc. IEEE DSP 2013*, 2013.
- [3] C. Breithaupt, et al., *IEEE Trans. ASLP* 2008.
- [4] Y. Murota, et al., *Proc. ICASSP*, pp. 7490-7494, 2014.
- [5] D. Kitamura, et al., *IEICE Trans. Fundamentals*, vol. E97-A, no. 5, pp. 1113-1118, 2014.