

**Abstract**: This paper presents a novel BSS approach that simultaneously performs an estimation of the number of sources, source separation based on the sparseness of speech, and

permutation alignment, based on Bayesian nonparametric approach.

## **1. Introduction**

### Blind source separation (BSS)

- Technique for separating sources only from microphone inputs
- Potential applications include hands-free teleconference system and automatic meeting transcription system

### Motivation and objective

- It is often difficult to pre-specify the exact number of all possible sources present in real environments.
- **e.g.)** In meeting situations, # of speakers can change during the meeting or loud, unexpected noise such as door slamming can occur.
- → When # of sources is unknown, we shall always consider underdetermined case (# of microphones < # of sources)
- One successful approach for underdetermined BSS involves utilizing "\*sparseness" of speech [1]
- \* T-F components of speech are near zero across most of T-F bins
- To exploit the sparseness of speech, mixing model must be represented in T-F domain
- → Permutation alignment problem needs to be solved

We propose BSS approach that simultaneously performs (1) estimation of # of sources, (2) source separation based on sparseness of speech, and (3) permutation alignment.

## 2. Mixing Model









# Blind separation of infinitely many sparse sources Hirokazu Kameoka<sup>1,2</sup>, Misa Sato<sup>1</sup>, Takuma Ono<sup>1</sup>, Nobutaka Ono<sup>3</sup>, Shigeki Sagayama<sup>1</sup> <sup>1</sup> The University of Tokyo, <sup>2</sup> NTT Communication Science Laboratories, <sup>3</sup> National Institute of Informatics





$$\hat{q}(A) = \prod_{k,\omega} \mathcal{N}_{\mathbb{C}}(\boldsymbol{a}_{k,\omega}; \boldsymbol{m}_{k,\omega}, \Gamma_{k,\omega})$$
$$\hat{q}(S) = \prod_{\omega,t} \mathcal{N}_{\mathbb{C}}(\hat{s}_{\omega,t}; \mu_{\omega,t}, \sigma_{\omega,t}^2)$$
$$\hat{q}(Z) = \prod \text{Discrete}(z_{\omega,t}; \boldsymbol{\phi}_{\omega,t})$$

 $^{\omega,t}$ 

 $\hat{q}(V) = \prod \text{Beta}(v_k; \gamma_{k,0}, \gamma_{k,1})$  $\hat{q}(C) = \prod \text{Discrete}(c_k; \boldsymbol{\psi}_k)$  $\hat{q}(\boldsymbol{\rho}) = \text{Dirichlet}(\boldsymbol{\rho}; \zeta_1, \dots, \zeta_I)$ 

## References

- [1] O. Yılmaz, S. Rickard, "Blind separation of speech mixtures via time-frequency masking," IEEE Trans. Signal Process., vol. 52, no. 7, pp. 1830–1847, 2004. [2] Y. Izumi, N. Ono, S. Sagayama, "Sparseness-based 2ch BSS using the EM algorithm in reverberant environment," in Proc. WASPAA, pp. 147–150, 2007. [3] J. Sethuraman, "A constructive definition of Dirichlet priors," *Statistica Sinica*, vol. 4, pp. 639–650, 1994. [4] T. S. Ferguson, "A Bayesian analysis of some nonparametric problems," Annals of Statistics, vol. 1, no. 2, pp. 209–230, 1973. [5] H. Sawada, S. Araki, S. Makino, "Underdetermined convolutive blind source separation via frequency bin-wise clustering and permutation alignment," *IEEE Trans. ASLP,* vol. 19, no. 3, pp. 516–527, 2010.
- [6] E. Vincent, R. Gribonval, C. F´evotte, "Performance measurement in blind audio source separation," IEEE Trans. ASLP, pp. 1462–1469, 2006.



