

# Multichannel Nonnegative Matrix Factorization in Convolutional Mixtures

## With Application to Blind Audio Source Separation

Alexey Ozerov<sup>1</sup> and Cédric Févotte<sup>2</sup>

<sup>1</sup>Institut Télécom, Télécom ParisTech, CNRS LTCI, Paris, France

<sup>2</sup>CNRS LTCI, Télécom ParisTech, Paris, France

e-mail: {ozero, fevotte}@telecom-paristech.fr

web: <http://perso.telecom-paristech.fr/~ozero/>

### Abstract

- ▶ Nonnegative Matrix Factorization (NMF) is usually used for **single-channel** audio signal power spectrogram decomposition.
- ▶ We propose a **multichannel NMF** framework in a general case of **convolutional mixtures** of sources.
- ▶ Possible applications: **source separation** (blind or supervised) and **information retrieval** from audio (e.g., music transcription).
- ▶ Here we apply multichannel NMF to different **stereo audio source separation** tasks, and obtain very **promising results**.

### Introduction

- ▶ Single-channel NMF for power spectrogram decomposition:

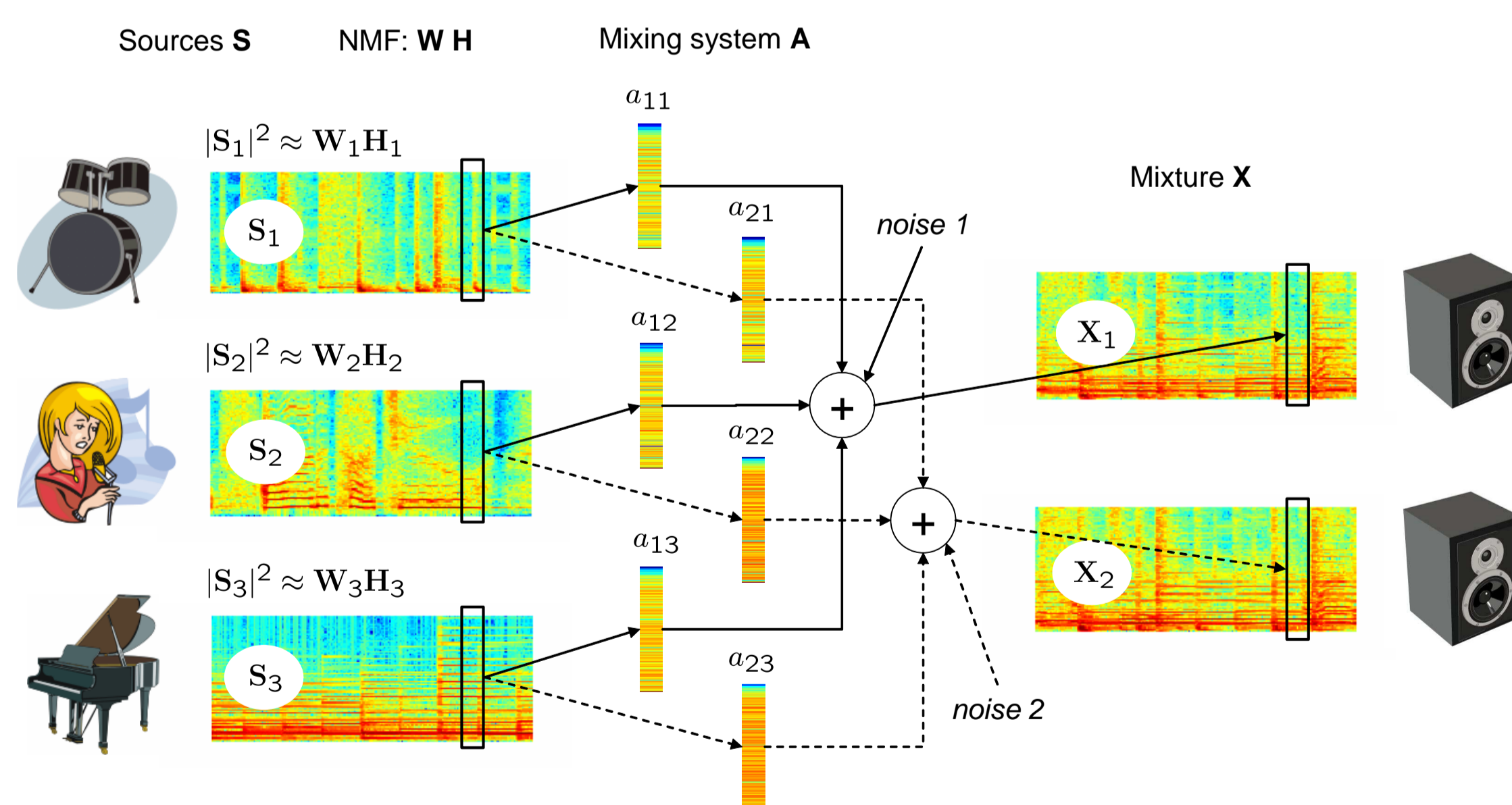
$$|\mathbf{X}|^2 \approx \mathbf{W}\mathbf{H}, \quad \mathbf{X} \in \mathbb{C}^{F \times N}, \quad \mathbf{W} \in \mathbb{R}_+^{F \times K}, \quad \mathbf{H} \in \mathbb{R}_+^{K \times N}$$

- ▶ Convolutional mixing equation (in STFT domain):

$$x_{i,fn} = \sum_{j=1}^J a_{ij,f} s_{j,fn} + b_{i,fn} \quad i = 1, \dots, I$$

mixture                      mix. coef    source            noise

- ▶ Multichannel NMF problem: **estimate A, W and H, given X.**
- ▶ Our approach is based on a **probabilistic NMF formulation**.



### Probabilistic NMF Model

- ▶ Source STFT is modeled as a sum of latent Gaussian components:

$$s_{j,fn} = \sum_{k=1}^{K_j} c_{j,k,fn} \quad \text{with} \quad c_{j,k,fn} \sim \mathcal{N}_c(0, w_{j,fk} h_{j,kn})$$

- ▶ **Maximum Likelihood (ML)** estimation of  $\mathbf{W}_j = [w_{j,fk}]_{f,k}$  and  $\mathbf{H}_j = [h_{j,kn}]_{k,n}$  given source STFT  $\mathbf{S}_j = [s_{j,fn}]_{f,n}$  is equivalent to NMF decomposition  $|\mathbf{S}_j|^2 \approx \mathbf{W}_j \mathbf{H}_j$  with **Itakura-Saito (IS) divergence** [1].

### Proposed Multichannel NMF Methods

#### Exact Likelihood Maximization with EM Algorithm

- ▶ Criterion  $(\theta = \{\mathbf{A}, \mathbf{W}, \mathbf{H}\}, \mathbf{x}_{fn} = [x_{1,fn}, \dots, x_{I,fn}]^T)$ :

$$C_1(\theta) = -\log p(\mathbf{X}|\theta) = -\sum_{fn} \log p(\mathbf{x}_{fn}|\theta).$$

- ▶ Expectation-Maximization (EM) Algorithm ( $\mathbf{C} = [c_{j,k,fn}]_{j,k,fn}$ ):

$$\begin{aligned} \text{E step:} \quad & Q(\theta|\theta^{(l)}) = \int \log p(\mathbf{X}, \mathbf{C}|\theta) p(\mathbf{C}|\mathbf{X}, \theta^{(l)}) d\mathbf{C}, \\ \text{M step:} \quad & \theta^{(l+1)} = \arg \max_{\theta} Q(\theta|\theta^{(l)}). \end{aligned}$$

- ▶  $\mathbf{A}^{(l+1)} \approx \mathbf{A}^{(l)}$  for small noise, thus a “simulated annealing” strategy is used.
- ▶ Related to other model-based methods (e.g., GMM-based [2]).

#### Individual Likelihoods Maximization with MU rules

- ▶ Criterion (IS divergence  $d_{IS}(x|y) = \frac{x}{y} - \log \frac{x}{y} - 1$ ):

$$\begin{aligned} C_2(\theta) &= -\sum_{i=1}^I \log p(\mathbf{X}_i|\theta) = -\sum_{i,fn} \log p(x_{i,fn}|\theta) \quad (\neq \log p(\mathbf{X}|\theta)) \\ &= \sum_{i,fn} d_{IS} \left( |x_{i,fn}|^2 \middle| \sum_j |a_{ij,f}|^2 \sum_{k=1}^{K_j} w_{j,fk} h_{j,kn} \right). \end{aligned}$$

- ▶ Multiplicative Update (MU) Rules:

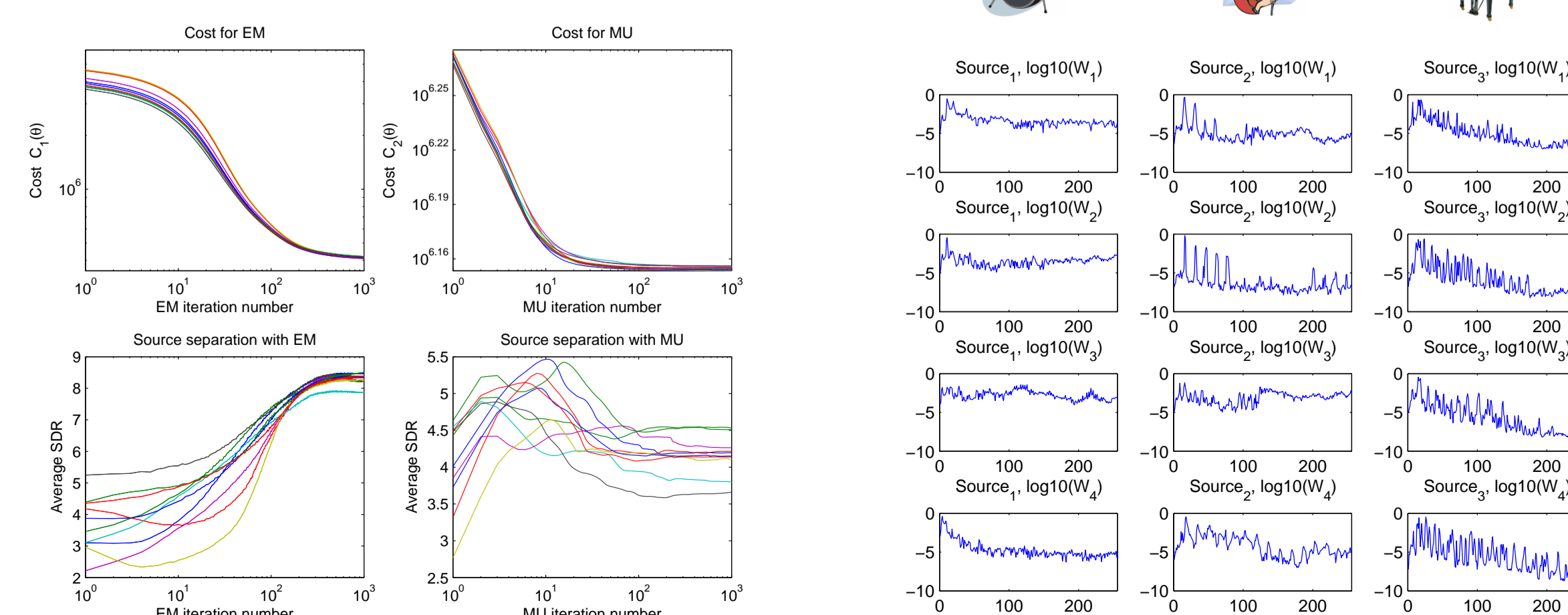
$$\theta_r \leftarrow \theta_r \frac{[\nabla_{\theta_r} C_2(\theta)]_-}{[\nabla_{\theta_r} C_2(\theta)]_+},$$

$$\nabla_{\theta_r} C_2(\theta) = [\nabla_{\theta_r} C_2(\theta)]_+ - [\nabla_{\theta_r} C_2(\theta)]_- \quad \text{and} \quad [\nabla_{\theta_r} C_2(\theta)]_+, [\nabla_{\theta_r} C_2(\theta)]_- \geq 0.$$

- ▶ Related to Nonnegative Tensor Factorization (NTF).

### Convergence and Decomposition

- ▶ Convergence vs. source separation performance (with perturbed oracle initializations).
- ▶ Multichannel NMF decomposition example (random init.).



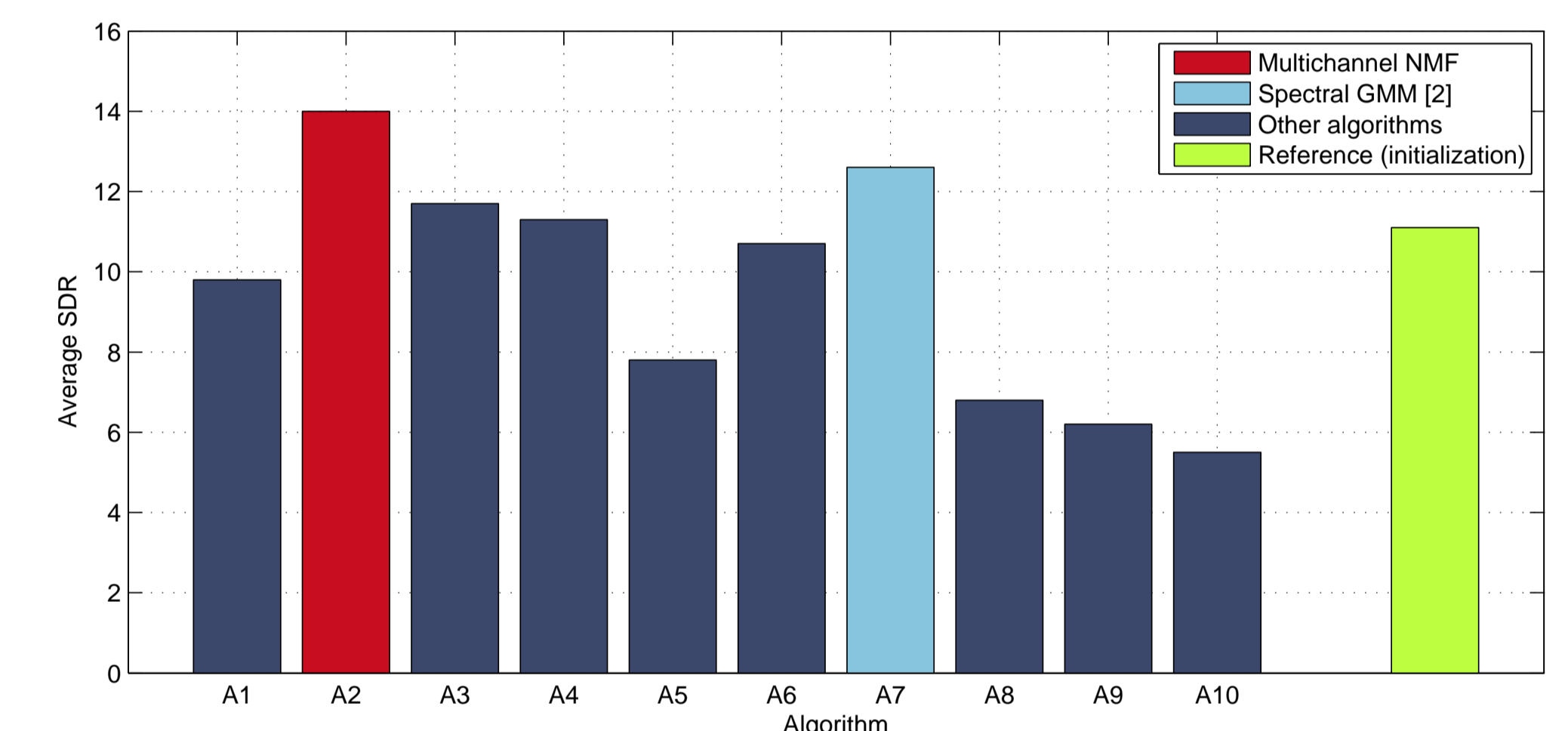
### Source Separation Results with EM

- ▶ Results for different mixture types (initialization = baseline).

| Mixture type                  | Music sources (3) |          | Speech sources (4) |          |
|-------------------------------|-------------------|----------|--------------------|----------|
|                               | Baseline          | Proposed | Baseline           | Proposed |
| Linear instantaneous          | 11.6              | 15.2     | 7.6                | 9.3      |
| Synthetic convolutional       | -0.8              | -0.6     | 3.5                | 4.5      |
| Live-recorded (convolutional) | 3.3               | 4.2      | 3.6                | 4.3      |

Table 1: Separation in terms of average Signal to Distortion Ratio (SDR) (dB).

- ▶ **Signal Separation Evaluation Campaign (SiSEC 2008) / ICA 2009.**
- ▶ “Under-determined speech and music (instantaneous) mixtures”:



### Conclusions

- ▶ **Strong points of the proposed approach:**
  - ▶ use of both **spectral** and **spatial** diversities for source separation,
  - ▶ **joint** and **blind** estimation of source and mixing models,
  - ▶ covers both **underdetermined** and **(over)determined noisy** cases,
  - ▶ source model frees us from convolutional BSS **permutation ambiguity**,
  - ▶ computational load **growing linearly** with number of components.
- ▶ **Weak point:** sensitive to the parameters initialization.

### References

- [1] C. Févotte, N. Bertin, and J.-L. Durrieu, “Nonnegative matrix factorization with the Itakura-Saito divergence. With application to music analysis,” *Neural Computation*, vol. 21, pp. 793–830, 2009.
- [2] S. Arberet, A. Ozerov, R. Gribonval, and F. Bimbot, “Blind spectral-GMM estimation for underdetermined instantaneous audio source separation,” in *ICA’09*, 2009.